

Package ‘forestinventory’

August 29, 2016

Type Package

Title Design-Based Global and Small-Area Estimations for Multiphase Forest Inventories

Version 0.1.0

Date 2016-07-06

Maintainer Andreas Hill <forestinventory@gmx.ch>

Description Extensive global and small-area estimation procedures for multiphase forest inventories under the design-based Monte-Carlo approach are provided. The implementation includes estimators for simple and cluster sampling published by Daniel Mandallaz in 2007 (<DOI:10.1201/9781584889779>), 2013 (<DOI:10.1139/cjfr-2012-0381>, <DOI:10.1139/cjfr-2013-0181>, <DOI:10.1139/cjfr-2013-0449>, <DOI:10.3929/ethz-a-009990020>) and 2016 (<DOI:10.3929/ethz-a-010579388>). It provides point estimates, their external- and design-based variances as well as confidence intervals. The procedures have also been optimized for the use of remote sensing data as auxiliary information.

License GPL (>= 2)

LazyData TRUE

Imports plyr (>= 1.8.3), stats, utils

RoxygenNote 5.0.1

NeedsCompilation no

Author Andreas Hill [aut, cre],
Alexander Massey [aut],
Daniel Mandallaz [ctb]

Repository CRAN

Date/Publication 2016-07-06 16:04:39

R topics documented:

confint	2
forestinventory	4

grisons	6
onephase	7
summary	9
threephase	10
twophase	17
zberg	23

Index	25
--------------	-----------

confint	<i>Calculates Confidence Intervals for Global and Small-Area Estimations</i>
---------	--

Description

Calculates Confidence Intervals for Global and Small-Area Estimations

Usage

```
## S3 method for class 'onephase'
confint(object, parm, level = 0.95,
        adjust.method = "none", ...)
```

```
## S3 method for class 'twophase'
confint(object, parm, level = 0.95,
        adjust.method = "none", ...)
```

```
## S3 method for class 'threephase'
confint(object, parm, level = 0.95,
        adjust.method = "none", ...)
```

Arguments

object	object of class onephase, twophase or threephase, containing estimation results of the respective estimation method.
parm	ignored.
level	the confidence level required.
adjust.method	correction method to obtain simultaneous confidence intervals for a set of estimates (thus restricted to objects of class "onephase", c("smallarea", "twophase") and c("smallarea", "threephase")). Available correction methods are c("none", "bonferroni"). Defaults to "none".
...	additional arguments, so far ignored.

Details

Depending on the estimation method specified, `confint()` computes confidence intervals as follows:

onephase:

Two-sided confidence intervals are computed based on the t -distribution with $n_2 - p$ *degrees of freedom*, where n_2 is the number of terrestrial data in the respective inventory domain.

twophase:

The calculation of the two-sided confidence intervals for *global* twophase estimates (objects of class `global`) are calculated based on the quantiles of the t -distribution with $n_2 - p$ *degrees of freedom*, where p is the number of parameters used in the regression model, and n_2 is the number of terrestrial observations (i.e. *local densities*) in the inventory domain.

The calculation of the two-sided confidence intervals for *smallarea* twophase estimates (objects of class `smallarea`) are calculated based on the quantiles of the t -distribution with $n_{2G} - 1$ *degrees of freedom*, where n_{2G} is the number of terrestrial observations (i.e. *local densities*) in the *smallarea*.

threephase:

The calculation of the two-sided confidence intervals for *global* threephase estimates (objects of class `global`) are calculated based on the quantiles of the t -distribution with $n_2 - p$ *degrees of freedom*, where p is the number of parameters used in the **full** regression model, and n_2 is the number of terrestrial observations (i.e. *local densities*) in the inventory domain (note: in notation used here n_0 , n_1 and n_2 correspond to the first, second and third phase sample sizes respectively).

The calculation of the two-sided confidence intervals for *smallarea* threephase estimates (objects of class `smallarea`) are calculated based on the quantiles of the t -distribution with $n_{2G} - 1$ *degrees of freedom*, where n_{2G} is the number of terrestrial observations (i.e. *local densities*) in the *smallarea*.

Value

`confint` returns a list of the following 3 components:

<code>ci</code>	a <code>data.frame</code> containing the columns: <ul style="list-style-type: none"> • <code>area</code> the domain, i.e. small area • <code>ci.lower.ext</code> the lower confidence limit based on the external variance • <code>ci.upper.ext</code> the upper confidence limit based on the external variance • <code>ci.lower.g</code> the lower confidence limit based on the g-weight variance • <code>ci.upper.g</code> the upper confidence limit based on the g-weight variance
<code>level</code>	the applied confidence level
<code>adjust.method</code>	the adjustment method applied to retrieve simultaneous confidence intervals

Note

In the special case of *synthetic* *smallarea* estimations, the two-sided confidence intervals are calculated based on the quantiles of the t -distribution with $n_2 - p$ *degrees of freedom*, i.e. based on the global sample size.

The confidence intervals for *synthetic* *smallarea* estimations do not account for the potential bias of a linear model that was fit in a large forest area and applied to a small area. Thus, the coverage

rates for confidence intervals produced by synthetic estimators may be less than the nominal level of confidence.

In case of cluster-sampling, `n2G` is the number of terrestrial clusters (a cluster constitutes the sample unit). This is automatically considered by `confint`.

The adjustment methods passed to `adjust.method` are designed to achieve *simultaneous* confidence intervals by correcting the confidence level given by `level`. The use of this option is recommended if a set of estimates contained in a `onophase`- or `smallarea`-object should be compared by their confidence intervals. It ensures that the percentage of confidence intervals containing the true value will correspond to the nominal confidence level.

References

Mandallaz, D. (2013). *Design-based properties of some small-area estimators in forest inventory with two-phase sampling*. Canadian Journal of Forest Research, 43(5), 441-449.

Mandallaz, D., Breschan, J., & Hill, A. (2013). *New regression estimators in forest inventories with two-phase sampling and partially exhaustive information: a design-based monte carlo approach with applications to small-area estimation*. Canadian Journal of Forest Research, 43(11), 1023-1031.

Mandallaz, D. (2013). *A three-phase sampling extension of the generalized regression estimator with partially exhaustive information*. Canadian Journal of Forest Research, 44(4), 383-388.

Benjamini, Y., and Hochberg, Y. (1995). *Controlling the false discovery rate: a practical and powerful approach to multiple testing*. Journal of the Royal Statistical Society Series B 57, 289-300.

Examples

```
## Calculate twophase estimations by extended pseudosynthetic estimator
# for 4 small areas ("A", "B", "C", "D") using the grisons-dataset:
sae.est <- twophase(formula = tvol ~ mean + stddev + max + q75,
                   data = grisons,
                   phase_id = list(phase.col = "phase_id_2p", terrgrid.id = 2),
                   small_area = list(sa.col = "smallarea",
                                     areas = c("A", "B", "C", "D"),
                                     unbiased = TRUE))

## calculate 95%-confidence intervals for each small area:
confint(sae.est)

## calculate simultaneous 95%-confidence intervals using 'bonferroni'-method:
confint(sae.est, adjust.method = "bonferroni")
```

Description

The package provides *global*- and *smallarea estimators* for *twophase* and *threephase* forest inventories under simple and cluster sampling, which have been developed by Daniel Mandallaz at ETH Zurich. The implemented methods have been published and applied in various studies (see References) and can be used for *double sampling for stratification*, *double sampling for regression* and *double sampling for regression within strata*.

Functions

The package provides three main functions to apply the various estimators for *twophase* and *threephase* forest inventories:

- [twophase](#) Function to apply global- and various smallarea estimation techniques for twophase inventories
- [threephase](#) Function to apply global- and various smallarea estimation techniques for three-phase inventories
- [onephase](#) Function to apply estimations for onephase inventories, mainly for comparison with two-and threephase

Motivation

The Motivation of writing this package was to provide an extensive and consistent collection of state-of-the-art *design-based* estimation techniques for forest inventories. It was especially designed to facilitate the application of the available estimators in forest practice as well as in scientifically related studies. The work on this package was also the trigger to complete the range of the already published estimators, especially in the framework of three-phase smallarea estimators.

Selected references

- Massey, A. F. (2015). *Multiphase estimation procedures for forest inventories under the design-based Monte Carlo approach* (Doctoral dissertation, Diss., ETH Zurich, Nr. 23025).
- Mandallaz, D. (2013). *Design-based properties of some small-area estimators in forest inventory with two-phase sampling*. Canadian Journal of Forest Research, 43(5), 441-449.
- Mandallaz, D., Breschan, J., & Hill, A. (2013). *New regression estimators in forest inventories with two-phase sampling and partially exhaustive information: a design-based monte carlo approach with applications to small-area estimation*. Canadian Journal of Forest Research, 43(11), 1023-1031.
- Mandallaz, D. (2013). *A three-phase sampling extension of the generalized regression estimator with partially exhaustive information*. Canadian Journal of Forest Research, 44(4), 383-388.

grisons

Data from a multiphase forest inventory in the canton of Grisons, Switzerland

Description

A dataset containing observations of 306 systematically arranged sample plots. Auxiliary information for all 306 plots is provided in the form of LiDAR canopy height metrics. For a systematic subsample of 67 out of the 306 plots, terrestrial information of the timber volume is provided from a terrestrial survey in the year 2007. Originally the inventory was carried out as a twophase inventory and has been artificially extended to a threephase inventory for demonstration purposes.

Usage

grisons

Format

data frame with 306 rows and 14 columns

Details

- `phase_id_2p` phase-membership of each observation for the twophase inventory. The large phase is indicated by 1, the terrestrial phase by 2.
- `phase_id_3p` phase-membership of each observation for the threephase inventory, i.e. the largest phase (0), the large phase (1) and terrestrial phase (2). *Note:* The threephase sample scheme was artificially created for demonstration purposes of the `threephase`-functions.
- `boundary_weights` proportion of analysis-window for auxiliary information lying within the forest.
- `mean` mean canopy height at the sample location based on the LiDAR canopy height model.
- `stddev` standard deviation of the LiDAR canopy height model at the sample location.
- `max` maximum value of the LiDAR canopy height model at the sample location.
- `q75` 75%-Quantile of the LiDAR canopy height model at the sample location.
- `smallarea` smallarea-indicator for each observation.
- `tv01` terrestrial timber volume from field survey. Use for `twophase`-inventory.
- `tv01.3p` terrestrial timber volume from field survey. Use for `threephase`-inventory.

Note

There are additional columns in `grisons` to demonstrate the function-behaviours for special cases which might occur in a forest inventory

- `phase_id_3p_nG0` one of the smallareas does not contain any terrestrial observation.
- `phase_id_3p_nG1` one of the smallareas does contain only a single terrestrial observation.

- `tvol.3p_nG0` Use as response variable to test `phase_id_3p_nG0` for [threephase-inventory](#).
- `tvol.3p_nG1` Use as response variable to test `phase_id_3p_nG1` for [threephase-inventory](#).

We leave testing these special cases to the user.

Source

The terrestrial data are kindly provided by the forest service of the canton grisons.

The dataset was created and used within the framework of the publications listed under *References*.

References

Mandallaz, D., Breschan, J., & Hill, A. (2013). *New regression estimators in forest inventories with two-phase sampling and partially exhaustive information: a design-based monte carlo approach with applications to small-area estimation*. Canadian Journal of Forest Research, 43(11), 1023-1031.

Hill, A., Breschan, J., & Mandallaz, D. (2014). *Accuracy assessment of timber volume maps using forest inventory data and LiDAR canopy height models*. Forests, 5(9), 2253-2275.

onephase

onephase

Description

`onephase` is used to calculate estimations exclusively based on terrestrial observations of a forest inventory (i.e. the *local densities*). The estimation method is available for simple and cluster-sampling and provides point estimates of the sample mean and their variances.

Usage

```
onephase(formula, data, phase_id = list(phase.col = NA, terrgrid.id = NA),
          cluster = NA, area = list(sa.col = NA, areas = NA))
```

Arguments

- | | |
|-----------------------|---|
| <code>formula</code> | an object of class " formula " that must be of the form $Y \sim 1$, where Y is the terrestrial response value of interest provided in data. |
| <code>data</code> | a data frame or vector containing the response value Y. Specifications are given under 'Details'. |
| <code>phase_id</code> | an object of class " list " containing two elements: <ul style="list-style-type: none"> • <code>phase.col</code>: the column name in data that specifies the phase membership of each observation • <code>terrgrid.id</code>: the indicator identifying the the terrestrial (a.k.a. "ground truth") phase for that column |

Note: Only has to be specified if data is of class `data.frame`.

cluster	Specifies the column name in data containing the cluster identification. Only used in case of cluster sampling.
area	(<i>Optional</i>) an object of class "list" containing two elements: <ul style="list-style-type: none"> • sa.col: the column name in data containing domain identification • areas: vector of desired domains for which the estimation should be computed. If estimations for multiple domains should be computed, the domains have to be defined within a character vector using c() Further details of the parameter-specifications are given under ' <i>Details</i> '.

Details

data can either be a vector only containing the observations of the response variable Y, or a data frame containing a column for the response variable and a column for the sample-grid indication that has to be further specified by argument phase_id. Additional *optional* columns include a cluster identification in case of cluster sampling, as well as a column that specifies a domain (e.g. a forest district) the respective terrestrial observation falls into. The latter allows to compute onephase-estimations for multiple domains at a time (see '*Examples*').

Value

onephase returns an object of class "onephase".

The functions `summary` and `confint` can be used to obtain a summary of the estimation results (point estimations, variances and sample sizes) and the confidence intervals for the respective point estimates.

An object of class "onephase" returns a list of the following components:

input	a list containing the function inputs
estimation	a data frame containing the following components: <ul style="list-style-type: none"> • area: the domain (only present if argument area has been used) • estimate: the point estimate • variance: the variance of the point estimate • n2: the terrestrial sample size
samplesizes	a named numeric vector giving the terrestrial samplesize

References

Mandallaz, D. (2007). *Sampling techniques for forest inventories*. Chapter 4. CRC Press.

Examples

```
# ----- non-cluster sampling-----#

## load grisons dataset:
data(grisons)

## 1) calculate onephase-estimation for entire dataset:
op <- onephase(formula = tvol~1 ,data = grisons,
```

```

                                phase_id =list(phase.col = "phase_id_2p",terrgrid.id = 2))
summary(op)
confint(op)

## 2) calculate onephase-estimation for given domains (areas) in dataset:
op.a <- onephase(formula = tvol~1,
                 data = grisons,
                 phase_id = list(phase.col = "phase_id_2p", terrgrid.id = 2),
                 area = list(sa.col = "smallarea", areas = c("A", "B")))
summary(op.a)
confint(op.a)

# ----- cluster sampling -----#

## load zurichberg dataset:
data(zberg)

## 1) calculate onephase-estimation for entire dataset:
op.clust <- onephase(formula = basal~1, data = zberg,
                    phase_id = list(phase.col = "phase_id_2p",terrgrid.id = 2),
                    cluster = "cluster")
summary(op.clust)
confint(op.clust)

## 2) calculate onephase-estimation for given areas in dataset:
op.clust.a <- onephase(formula = basal~1,
                      data = zberg,
                      phase_id = list(phase.col = "phase_id_2p", terrgrid.id = 2),
                      cluster = "cluster",
                      area = list(sa.col = "ismallg23", areas = c("2", "3")))
summary(op.clust.a)
confint(op.clust.a)

```

summary

Summarizing Global and Small-Area Estimation Results

Description

Summarizing Global and Small-Area Estimation Results

Usage

```
## S3 method for class 'onephase'
summary(object, coefs = FALSE, ...)
```

```
## S3 method for class 'twophase'
summary(object, coefs = FALSE, ...)
```

```
## S3 method for class 'threephase'
summary(object, coefs = FALSE, ...)
```

Arguments

object	object of class <code>onephase</code> , <code>twophase</code> or <code>threephase</code> , containing estimation results of the respective estimation method.
coefs	of type <code>"logical"</code> . If set to <code>TRUE</code> , also gives the regression coefficients of <code>twophase</code> and <code>threephase</code> estimations. Defaults to <code>FALSE</code> .
...	additional arguments, so far ignored.

threephase	<i>threephase</i>
------------	-------------------

Description

`threephase` is used to calculate estimations based on triple sampling under the *model-assisted Monte Carlo approach*. A *first phase* of auxiliary information (e.g. taken from remote sensing data) is used to generate model predictions based on multiple linear regression using the method of ordinary least squares. A subsample of the first phase comprises a *second phase* which contains further auxiliary information that produces another set of model predictions. A further subsample produces a *third final phase* based on terrestrial observations (i.e. the *local densities* of the ground truth) and is used to correct for bias in the design-based sense. The estimation method is available for *simple* and *cluster sampling* and includes the special case where the first phase is based on an *exhaustive* sample (i.e. a census). *Small-area applications* are supported for synthetic estimation as well as two varieties of bias-corrected estimators: the traditional small-area estimator and an asymptotically equivalent version derived under Mandallaz's extended model approach.

Usage

```
threephase(formula.s0, formula.s1, data, phase_id, cluster = NA,
  small_area = list(sa.col = NA, areas = NA, unbiased = TRUE),
  boundary_weights = NA, exhaustive = NA, progressbar = FALSE,
  psmall = FALSE)
```

Arguments

<code>formula.s0</code>	an object of class <code>"formula"</code> as would be used in the function <code>lm</code> that contains a reduced set of auxiliary variables available for all first phase plots
<code>formula.s1</code>	an object of class <code>"formula"</code> as would be used in the function <code>lm</code> that contains the predictors from <code>formula.s0</code> as well as further ancilliary predictors available for all second phase plots (i.e. <code>formula.s0</code> is nested in <code>formula.s1</code>)
<code>data</code>	a data frame containing all variables contained in <code>formula</code> and a column indexing phase membership. Additional columns designating small-area membership, cluster ID and boundary weights should also be contained in the data frame if they are requested in the function.
<code>phase_id</code>	an object of class <code>"list"</code> containing three elements: <ul style="list-style-type: none"> • <code>phase.col</code>: the column name in <code>data</code> that specifies the phase membership of each observation

	<ul style="list-style-type: none"> • <code>s1.id</code>: the indicator identifying the "second phase only" plots for that column • <code>terrgrid.id</code>: the indicator identifying the terrestrial (a.k.a. "ground truth") phase for that column
<code>cluster</code>	(<i>Optional</i>) Specifies the column name in data containing the cluster ID. Only used in case of cluster sampling.
<code>small_area</code>	(<i>Optional</i>) a list that if containing three elements: <ul style="list-style-type: none"> • <code>sa.col</code>: the column name in data containing domain identification • <code>areas</code>: vector of desired small-area domain identifiers • <code>unbiased</code>: an object of type "logical" that when FALSE designates that the estimator is allowed to be biased (i.e. the synthetic estimator) and when TRUE forces it to be design-unbiased. See '<i>Details</i>'. <p>Note: If <code>small_area</code> is left unchanged then <code>twophase</code> defaults to global estimation.</p>
<code>boundary_weights</code>	(<i>Optional</i>) Specifies the column name in data containing the weights for boundary adjustment. See ' <i>Details</i> '
<code>exhaustive</code>	(<i>Optional</i>) For global estimation, a vector of true auxiliary means corresponding to an exhaustive first phase. The vector must be input in the same order that <code>lm</code> processes a formula object and include the intercept term. For small area estimation, <code>exhaustive</code> is a <code>data.frame</code> containing column names (<code>colnames</code>) for every variable appearing in the parameter formula including the variable "Intercept". Rownames (<code>rownames</code>) have to be used and must correspond to the names of the small areas. See ' <i>Details</i> '.
<code>progressbar</code>	(<i>Optional</i>) an object a type "logical" that when TRUE prints the progress of the calculation in the console (recommended for large amount of small areas). Defaults to FALSE.
<code>psmall</code>	(<i>Optional</i>) an object a type "logical" used for small area estimations that only works when <code>unbiased</code> in the parameter <code>small_area</code> is set to TRUE. See ' <i>Details</i> '.

Details

`s1.id` identifies "second phase only" plots because the terrestrial phase is known to be part of the second phase by the construction of the subsampling.

If estimations for multiple small-area domains should be computed, the domains have to be defined within a character vector using `c()`. Using `small_area(..., unbiased=FALSE)` calculates design-based estimates with the synthetic estimator and may be design-biased if the model is biased in that small area. The default, `small_area(..., unbiased=TRUE)`, allows for a residual correction by one of two asymptotically equivalent methods to create design-unbiased estimates:

- Mandallaz's extended model approach calculates the residual correction by extending the model formula with an indicator variable in the small area. It is the default method `psmall=FALSE`.
- the traditional small area estimator calculates the residual correction by taking the synthetic estimator and adding the mean residual observed in the small area. It is activated when `psmall=TRUE`.

Missing values (NA) in the auxiliary variables (i.e. at least one auxiliary variable cannot be observed at an inventory location) are automatically removed from the dataset *before* the estimations are computed. Note that missingness in the auxiliary variables is only allowed if we assume that they are *missing at random*, since the unbiasedness of the estimates is based on the sampling design.

The boundary weight adjustment is pertinent for auxiliary information derived from remote sensing and is equal to the percentage of forested area (e.g. as defined by a forest mask) in the interpretation area.

Exhaustive estimation refers to when the true means of certain auxiliary variables are known and an exhaustive first phase (i.e. a census). For global estimation, the vector must be input in the same order that `lm` processes a `formula` object including the intercept term whose true mean will always be one. For small area estimation, `exhaustive` is a `data.frame` containing column names for every variable appearing in the parameter `formula` including the variable "Intercept". The observations of the `data.frame` must represent the true auxiliary means in the same order as was presented in areas from the parameter `small_area`. See '*Examples*'.

Value

`threephase` returns an object of class "threephase".

An object of class "threephase" returns a list of the following components:

<code>input</code>	a list containing the function's inputs
<code>estimation</code>	a data frame containing the following components: <ul style="list-style-type: none"> • <code>area</code>: the domain (only present if argument <code>areas</code> has been used) • <code>estimate</code>: the point estimate • <code>ext_variance</code>: the external variance of the point estimate that doesn't account for fitting the model from the current inventory • <code>g_variance</code>: the internal (g-weight) variance that accounts for fitting the model from the current inventory • <code>n0</code> the first phase sample size of plots • <code>n1</code> the second phase sample size of plots • <code>n2</code> the third phase (i.e. terrestrial) sample size of plots • <code>n0G</code> the first phase sample size in the small area • <code>n1G</code> the second phase sample size in the small area • <code>n2G</code> the third phase (i.e. terrestrial) sample size in the small area • <code>r_squared_reduced</code> the R-squared of the linear model based on <code>formula.s0</code> (i.e. the reduced model) • <code>r_squared_full</code> the R-squared of the linear model based on <code>formula.s1</code> (i.e. the full model)
<code>samplesizes</code>	a <code>data.frame</code> summarizing all <code>samplesizes</code> : in case of cluster sampling both, the number of individual plots and the number of clusters is reported.
<code>coefficients</code>	the coefficients of the two linear models: <ul style="list-style-type: none"> • <code>alpha</code>: the reduced model coefficients • <code>beta</code>: the full model coefficients
<code>cov_alpha_s2</code>	the design-based covariance matrix of the reduced model coefficients

cov_beta_s2	the design-based covariance matrix of the full model coefficients
Z_bar_1_s0	the estimated auxiliary means of formula.s0 based on the first phase. If the first phase is exhaustive, these are the true auxiliary means specified in the input-argument exhaustive.
Z1_bar_s1	the estimated auxiliary means of formula.s0 based on the second phase
Z_bar_s1	the estimated auxiliary means of formula.s1 based on the second phase
cov_Z_bar_1_s0	the covariance matrix for Z_bar_1_s0
resid_reduced	the reduced model residuals at either the plot level or cluster level depending on the call
resid_full	the full model residuals at either the plot level or cluster level depending on the call
warn.messages	logical indicating if warning messages were issued

Note

In the special case of cluster sampling, the reported sample sizes in estimation are the number of clusters. The `samplesize`-object also provides the respective number of single plot units for cluster sampling. The reported `r.squared_reduced` and `r.squared_full` describe the model fit of the applied linear regression models (i.e. on *plot-level*, not on *cluster level*).

References

- Mandallaz, D., Breschan, J., & Hill, A. (2013). *New regression estimators in forest inventories with two-phase sampling and partially exhaustive information: a design-based monte carlo approach with applications to small-area estimation*. Canadian Journal of Forest Research, 43(11), 1023-1031.
- Mandallaz, D. (2014). *A three-phase sampling extension of the generalized regression estimator with partially exhaustive information*. Can. J. For. Res. 44: 383-388
- Massey, A. and Mandallaz, D. and Lanz, A. (2014). *Integrating remote sensing and past inventory data under the new annual design of the Swiss National Forest Inventory using three-phase design-based regression estimation*. Can. J. For. Res. 44(10): 1177-1186
- Mandallaz, D. (2013). *Regression estimators in forest inventories with three-phase sampling and two multivariate components of auxiliary information*. ETH Zurich, Department of Environmental Systems Science, Tech. rep. Available from <http://e-collection.library.ethz.ch>.

Examples

```
## load datasets:
data(grisons)
data(zberg)

## define regression models for simple and cluster sampling:
formula.s0 <- tvol ~ mean # reduced model:
formula.s1 <- tvol ~ mean + stddev + max + q75 # full model
formula.clust.s0 <- basal ~ stade
```

```

formula.clust.s1 <- basal ~ stade + couver + melange

# -----#
# ----- GLOBAL ESTIMATION -----#

#----
## 1) -- Design-based estimation with non-exhaustive auxiliary information
#----

# 1.1) non-cluster-sampling (see eqns. [11], [14] and [16] in Mandallaz 2014):
summary(threephase(formula.s0, formula.s1, data = grisons,
                   phase_id = list(phase.col = "phase_id_3p", s1.id=1, terrgrid.id = 2)))

# 1.2) cluster-sampling (see eqns. [49] and [50] in Mandallaz 2013):
summary(threephase(formula.clust.s0, formula.clust.s1, data = zberg,
                   phase_id = list(phase.col="phase_id_3p", s1.id = 1, terrgrid.id = 2),
                   cluster = "cluster"))

# 1.3) example for boundary weight adjustment (non-cluster example):
summary(threephase(formula.s0, formula.s1, data = grisons,
                   phase_id = list(phase.col="phase_id_3p", s1.id = 1, terrgrid.id = 2),
                   boundary_weights = "boundary_weights"))

#----
## 2) -- Design-based estimation with exhaustive auxiliary information
#----

# 2.1) non-cluster-sampling (see eqns. [7], [9] and [10] in Mandallaz 2014):
summary(threephase(formula.s0, formula.s1, data = grisons,
                   phase_id = list(phase.col = "phase_id_3p", s1.id = 1, terrgrid.id = 2),
                   exhaustive = c(1,11.39)))

# 2.2) cluster-sampling:
summary(threephase(formula.clust.s0, formula.clust.s1, data = zberg,
                   phase_id = list(phase.col = "phase_id_3p", s1.id = 1, terrgrid.id = 2),
                   cluster = "cluster", exhaustive = c(1, 0.10, 0.7, 0.10)))

# -----#
# ----- SMALL AREA ESTIMATION -----#

#----
## 1) -- Design-based estimation with non-exhaustive auxiliary information
#----

# 1.1) Mandallaz's extended pseudo small area estimator:
summary(threephase(formula.s0,
                   formula.s1,
                   data = grisons,

```

```

phase_id = list(phase.col = "phase_id_3p", s1.id = 1, terrgrid.id = 2),
small_area=list(sa.col = "smallarea", areas = c("A", "B", "C", "D"),
unbiased = TRUE)))

summary(threephase(formula.clust.s0,
formula.clust.s1,
data = zberg,
phase_id = list(phase.col = "phase_id_3p", s1.id = 1, terrgrid.id = 2),
cluster = "cluster",
small_area = list(sa.col = "ismallold", areas = c("1"), unbiased = TRUE)))

# 1.2) pseudo small area estimator:
summary(threephase(formula.s0,
formula.s1,
data = grisons,
phase_id = list(phase.col = "phase_id_3p", s1.id = 1, terrgrid.id = 2),
small_area = list(sa.col = "smallarea", areas = c("A", "B", "C", "D"),
unbiased = TRUE),
psmall = TRUE))

summary(threephase(formula.clust.s0,
formula.clust.s1,
data=zberg,
phase_id=list(phase.col="phase_id_3p", s1.id=1, terrgrid.id=2),
cluster="cluster",
small_area=list(sa.col="ismallold", areas=c("1"), unbiased=TRUE),
psmall = TRUE))

# 1.3) pseudosynthetic small area estimator:
summary(threephase(formula.s0 = tvol ~ mean,
formula.s1 = tvol ~ mean + stddev + max + q75,
data = grisons,
phase_id = list(phase.col = "phase_id_3p", s1.id = 1, terrgrid.id = 2),
small_area = list(sa.col = "smallarea", areas = c("A", "B", "C", "D"),
unbiased = FALSE)))

summary(threephase(formula.clust.s0,
formula.clust.s1,
data = zberg,
phase_id = list(phase.col = "phase_id_3p", s1.id = 1, terrgrid.id = 2),
cluster = "cluster",
small_area = list(sa.col = "ismallold", areas = c("1"), unbiased = FALSE)))

#----
## 2) -- Design-based estimation with exhaustive auxiliary information
#----

# true auxiliary mean for variable "mean" taken from Mandallaz et al. (2013):
truemeans.G <- data.frame(Intercept = rep(1, 4),
mean = c(12.85, 12.21, 9.33, 10.45))

```



```

        exhaustive = truemmeans.G))

summary(threephase(formula.clust.s0,
                   formula.clust.s1,
                   data = zberg,
                   phase_id = list(phase.col = "phase_id_3p", s1.id = 1, terrgrid.id = 2),
                   cluster = "cluster",
                   small_area = list(sa.col = "ismallold", areas = c("1"), unbiased = FALSE),
                   exhaustive = truemmeans.G.clust))

```

twophase

twophase

Description

twophase is used to calculate estimations based on double sampling under the *model-assisted Monte Carlo approach*. A *first phase* of auxiliary information (e.g. taken from remote sensing data) is used to generate model predictions based on multiple linear regression using the method of ordinary least squares. A subsample of the first phase comprises the *second phase* which contains terrestrial observations (i.e. the *local densities* of the ground truth) that is used to correct for bias in the design-based sense. The estimation method is available for *simple* and *cluster sampling* and includes the special case where the first phase is based on an *exhaustive* sample (i.e. a census). *Small-area applications* are supported for synthetic estimation as well as two varieties of bias-corrected estimators: the traditional small-area estimator and an asymptotically equivalent version derived under Mandallaz's extended model approach.

Usage

```

twophase(formula, data, phase_id, cluster = NA, small_area = list(sa.col =
  NA, areas = NA, unbiased = TRUE), boundary_weights = NA, exhaustive = NA,
  progressbar = FALSE, psmall = FALSE)

```

Arguments

formula	an object of class " formula " as would be used in the function lm
data	a data frame containing all variables contained in formula and a column indexing phase membership. Additional columns designating small-area membership, cluster ID and boundary weights should also be contained in the data frame if they are requested in the function.
phase_id	an object of class " list " containing two elements: <ul style="list-style-type: none"> • <code>phase.col</code>: the column name in data that specifies the phase membership of each observation • <code>terrgrid.id</code>: the indicator identifying the terrestrial (a.k.a. "ground truth") phase for that column
cluster	(<i>Optional</i>) Specifies the column name in data containing the cluster ID. Only used in case of cluster sampling.

small_area	<p>(<i>Optional</i>) a list that if containing three elements:</p> <ul style="list-style-type: none"> • sa.col: the column name in data containing domain identification • areas: vector of desired small-area domain identifiers • unbiased: an object of type "logical" that when FALSE designates that the estimator is allowed to be biased (i.e. the synthetic estimator) and when TRUE forces it to be design-unbiased. See '<i>Details</i>'. <p>Note: If small_area is left unchanged then twophase defaults to global estimation.</p>
boundary_weights	<p>(<i>Optional</i>) Specifies the column name in data containing the weights for boundary adjustment. See '<i>Details</i>'</p>
exhaustive	<p>(<i>Optional</i>) For global estimation, a vector of true auxiliary means corresponding to an exhaustive first phase. The vector must be input in the same order that lm processes a formula object and include the intercept term. For small area estimation, exhaustive is a data.frame containing column names (colnames) for every variable appearing in the parameter formula including the variable "Intercept". Rownames (row.names) have to be used and must correspond to the names of the small areas. See '<i>Details</i>'.</p>
progressbar	<p>(<i>Optional</i>) an object a type "logical" that when TRUE prints the progress of the calculation in the console (recommended for large amount of small areas). Defaults to FALSE.</p>
psmall	<p>(<i>Optional</i>) an object a type "logical" used for small area estimations that only works when unbiased in the parameter small_area is set to TRUE. See '<i>Details</i>'.</p>

Details

If estimations for multiple small-area domains should be computed, the domains have to be defined within a character vector using `c()`. Using `small_area(..., unbiased=FALSE)` calculates design-based estimates with the synthetic estimator and may be design-biased if the model is biased in that small area. The default, `small_area(..., unbiased=TRUE)`, allows for a residual correction by one of two asymptotically equivalent methods to create design-unbiased estimates:

- Mandallaz's extended model approach calculates the residual correction by extending the model formula with an indicator variable in the small area. It is the default method `psmall=FALSE`.
- the traditional small area estimator calculates the residual correction by taking the synthetic estimator and adding the mean residual observed in the small area. It is activated when `psmall=TRUE`.

Missing values (NA) in the auxiliary variables (i.e. at least one auxiliary variable cannot be observed at an inventory location) are automatically removed from the dataset *before* the estimations are computed. Note that missingness in the auxiliary variables is only allowed if we assume that they are *missing at random*, since the unbiasedness of the estimates is based on the sampling design.

The boundary weight adjustment is pertinent for auxiliary information derived from remote sensing and is equal to the percentage of forested area (e.g. as defined by a forest mask) in the interpretation area.

Exhaustive estimation refers to when the true means of certain auxiliary variables are known and an exhaustive first phase (i.e. a census). For global estimation, the vector must be input in the same order that `lm` processes a formula object including the intercept term whose true mean will always be one. For small area estimation, `exhaustive` is a `data.frame` containing column names for every variable appearing in the parameter formula including the variable "Intercept". The observations of the `data.frame` must represent the true auxiliary means in the same order as was presented in areas from the parameter `small_area`. See '*Examples*'.

Value

`twophase` returns an object of class "`twophase`".

An object of class "`twophase`" returns a list of the following components:

<code>input</code>	a list containing the function's inputs
<code>estimation</code>	a data frame containing the following components: <ul style="list-style-type: none"> • <code>area</code>: the domain (only present if argument <code>areas</code> has been used) • <code>estimate</code>: the point estimate • <code>ext_variance</code>: the external variance of the point estimate that doesn't account for fitting the model from the current inventory • <code>g_variance</code>: the internal (g-weight) variance that accounts for fitting the model from the current inventory • <code>n1</code> the first phase sample size of plots • <code>n2</code> the second phase (i.e. terrestrial) sample size of plots • <code>n1G</code> the first phase sample size in the small area • <code>n2G</code> the second phase (i.e. terrestrial) sample size in the small area • <code>r_squared</code> the R squared of the linear model
<code>samplesizes</code>	a <code>data.frame</code> summarizing all samplesizes: in case of cluster sampling both, the number of individual plots and the number of clusters is reported.
<code>coefficients</code>	the linear model coefficients
<code>cov_coef</code>	the design-based covariance matrix of the model coefficients
<code>Z_bar_1G</code>	the estimated auxiliary means of formula based on the first phase. If the first phase is exhaustive, these are the true auxiliary means specified in the input-argument <code>exhaustive</code> .
<code>cov_Z_bar_1G</code>	the covariance matrix of <code>Z_bar_1G</code>
<code>Rc_x_hat_G</code>	the small-area residuals at either the plot level or cluster level depending on the call
<code>Rc_x_hat</code>	the residuals at either the plot level or cluster level depending on the call
<code>Yx_s2G</code>	the local densities in the small area
<code>Mx_s2G</code>	the cluster weights in the small area
<code>mean_Rc_x_hat_G</code>	the mean residual (weighted mean in the case of cluster sampling) in the small area
<code>mean_Rc_x_hat</code>	the mean residual (weighted mean in the case of cluster sampling)
<code>warn.messages</code>	logical indicating if warning messages were issued

Note

In the special case of cluster sampling, the reported sample sizes in estimation are the number of clusters. The `samplesize`-object also provides the respective number of single plot units for cluster sampling. The reported `r.squared` describe the model fit of the applied linear regression model (i.e. on *plot-level*, not on *cluster level*).

References

- Mandallaz, D. (2007). *Sampling techniques for forest inventories*. Chapter 4. CRC Press.
- Mandallaz, D. (2013). *Design-based properties of some small-area estimators in forest inventory with two-phase sampling*. Can. J. For. Res. 43: 441-449
- Mandallaz, D. and Hill, A. and Massey, A. (2016). *Design-based properties of some small-area estimators in forest inventory with two-phase sampling*. ETH Zurich, Department of Environmental Systems Science, Tech. rep. Available from <http://e-collection.library.ethz.ch>.

Examples

```
## load datasets:
data(grisons)
data(zberg)

# -----#
# ----- GLOBAL ESTIMATION -----#

#----
## 1) -- Design-based estimation with non-exhaustive auxiliary information
#----

# 1.1) non-cluster-sampling:
summary(twophase(formula = tvol ~ mean + stddev + max + q75,
                 data = grisons,
                 phase_id = list(phase.col = "phase_id_2p", terrgrid.id = 2)))

# 1.2) cluster-sampling (see eqns. [57] and [58] in Mandallaz, Hill, Massey 2016):
summary(twophase(formula = basal ~ stade + couver + melange,
                 data = zberg,
                 phase_id = list(phase.col = "phase_id_2p", terrgrid.id = 2),
                 cluster = "cluster"))

# 1.3) example for boundary weight adjustment (non-cluster example):
summary(twophase(formula=tvol ~ mean + stddev + max + q75,
                 data=grisons,
                 phase_id=list(phase.col = "phase_id_2p", terrgrid.id = 2),
                 boundary_weights = "boundary_weights"))

#----
## 2) -- Design-based estimation with exhaustive auxiliary information
#----

# establish order for vector of true auxiliary means:
```

```

colnames(lm(formula = tvol ~ mean + stddev + max + q75, data = grisons, x = TRUE)$x)
true.means <- c(1, 11.39, 8.84, 32.68, 18.03)

# 2.1) non-cluster-sampling:
summary(twophase(formula = tvol ~ mean + stddev + max + q75,
                 data = grisons,
                 phase_id = list(phase.col = "phase_id_2p", terrgrid.id = 2),
                 exhaustive = true.means))

# 2.2) cluster-sampling:
summary(twophase(formula = stem ~ stade + couver + melange,
                 data = zberg,
                 phase_id = list(phase.col = "phase_id_2p", terrgrid.id = 2),
                 cluster = "cluster",
                 exhaustive = c(1, 0.10, 0.7, 0.10, 0.6, 0.8)))

# -----#
# ----- SMALL AREA ESTIMATION -----#

#----
## 1) -- Design-based estimation with non-exhaustive auxiliary information
#----

# 1.1) Mandallaz's extended pseudo small area estimator (see eqns. [35] and [36] in Mandallaz 2013):
summary(twophase(formula = tvol ~ mean + stddev + max + q75, data = grisons,
                 phase_id = list(phase.col = "phase_id_2p", terrgrid.id = 2),
                 small_area = list(sa.col = "smallarea", areas = c("A", "B", "C", "D"),
                                   unbiased = TRUE)))

summary(twophase(formula = basal ~ stade + couver + melange, data=zberg,
                 phase_id = list(phase.col = "phase_id_2p", terrgrid.id = 2),
                 cluster = "cluster",
                 small_area = list(sa.col = "ismallg23", areas = c("2", "3"),
                                   unbiased = TRUE)))

# 1.2) pseudo small area estimator (see eqns. [25] and [26] in Mandallaz 2013):
summary(twophase(formula = tvol ~ mean + stddev + max + q75, data = grisons,
                 phase_id = list(phase.col = "phase_id_2p", terrgrid.id = 2),
                 small_area = list(sa.col = "smallarea", areas = c("A", "B"),
                                   unbiased = TRUE),
                 psmall = TRUE))

summary(twophase(formula = basal ~ stade + couver + melange, data=zberg,
                 phase_id = list(phase.col = "phase_id_2p", terrgrid.id = 2),
                 cluster = "cluster",
                 small_area = list(sa.col = "ismallg23", areas = c("2", "3"),
                                   unbiased = TRUE),
                 psmall = TRUE))

# 1.3) pseudosynthetic small area estimator (see eqns. [35] and [36] in Mandallaz 2013):

```

```

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

summary(twophase(formula = basal ~ stade + couver + melange, data=zberg,
  phase_id = list(phase.col = "phase_id_2p", terrgrid.id = 2),
  cluster = "cluster",
  small_area = list(sa.col = "ismallg23", areas = c("2", "3"),
    unbiased = FALSE)))

#----
## 2) -- Design-based estimation with exhaustive auxiliary information
#----

# establish order for vector of true auxiliary means:
colnames(lm(formula = tvol ~ mean + stddev + max + q75, data = grisons, x = TRUE)$x)
colnames(lm(formula = basal ~ stade + couver + melange, data = zberg, x = TRUE)$x)

# true auxiliary means taken from Mandallaz et al. (2013):
truemeans.G <- data.frame(Intercept = rep(1, 4),
  mean = c(12.85, 12.21, 9.33, 10.45),
  stddev = c(9.31, 9.47, 7.90, 8.36),
  max = c(34.92, 35.36, 28.81, 30.22),
  q75 = c(19.77, 19.16, 15.40, 16.91))
rownames(truemeans.G) <- c("A", "B", "C", "D")

# true auxiliary means taken from Mandallaz (1991):
truemeans.G.clust <- data.frame(Intercept = 1,
  stade400 = 0.175,
  stade500 = 0.429,
  stade600 = 0.321,
  couver2 = 0.791,
  melange2 = 0.809)
rownames(truemeans.G.clust) <- c("1")

# 2.1) Mandallaz's extended small area estimator (see eqns. [31] and [33] in Mandallaz 2013):
summary(twophase(formula = tvol ~ mean + stddev + max + q75, data = grisons,
  phase_id = list(phase.col = "phase_id_2p", terrgrid.id = 2),
  small_area = list(sa.col = "smallarea", areas = c("A", "B"),
    unbiased = TRUE),
  exhaustive = truemeans.G))

summary(twophase(formula = basal ~ stade + couver + melange, data=zberg,
  phase_id = list(phase.col = "phase_id_2p", terrgrid.id = 2),
  cluster = "cluster",
  small_area = list(sa.col = "ismalloid", areas = c("1"),
    unbiased = TRUE),
  exhaustive = truemeans.G.clust))

```

```
# 2.2) small area estimator (see eqns. [20] and [21] in Mandallaz 2013):
summary(twophase(formula = tvol ~ mean + stddev + max + q75, data = grisons,
  phase_id = list(phase.col = "phase_id_2p", terrgrid.id = 2),
  small_area = list(sa.col = "smallarea", areas = c("A"),
    unbiased = TRUE),
  exhaustive = truemmeans.G, psmall = TRUE))

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

# 2.3) synthetic small area estimator (see eqns. [18] and [19] in Mandallaz 2013):
summary(twophase(formula=tvoll ~ mean + stddev + max + q75, data=grisons,
  phase_id=list(phase.col = "phase_id_2p", terrgrid.id = 2),
  small_area=list(sa.col = "smallarea", areas = c("A", "B"),
    unbiased = FALSE),
  exhaustive = truemmeans.G))

summary(twophase(formula = basal ~ stade + couver + melange, data = zberg,
  phase_id = list(phase.col = "phase_id_2p", terrgrid.id = 2),
  cluster = "cluster",
  small_area = list(sa.col = "ismallold", areas = c("1"),
    unbiased = FALSE),
  exhaustive = truemmeans.G.clust))
```

zberg	<i>Data from a multiphase forest inventory at the zurichberg (zurich), switzerland</i>
-------	--

Description

A dataset from 1991 containing 1203 sample plots observations from a forest inventory using cluster-sampling. The large phase comprises 298 clusters. Terrestrial information of the stem number as well as the basal area is available for a systematic subsample of 73 clusters. Auxiliary information at all 2103 sample plots were derived by stand maps. Originally the inventory was carried out as a twophase inventory and has been artificially extended to a threephase inventory for demonstration purposes.

Usage

```
zberg
```

Format

data frame with 1203 rows and 12 columns

Details

- cluster cluster identification. Maximum number of sample plots per cluster is 5.
- phase_id_2p phase-membership of each observation for the twophase inventory. The first phase is indicated by 1, the second (i.e. terrestrial) phase by 2.
- phase_id_3p the phase-membership of each observation for the threephase inventory, i.e. the first phase (0), the second phase (1) and third (terrestrial) phase (2). *Note:* The threephase sample scheme was artificially created for demsontration purposes of the [threephase](#)-functions.
- stade development stage at sample plot location based on the stand map. Categorical variable of class factor with 4 levels.
- melange degree of mixture at sample plot location based on the stand map. Categorical variable of class factor with 2 levels.
- couver crown-coverage at sample plot location based on the stand map. Categorical variable of class factor with 2 levels.
- stem stem number dervied at field survey.
- basal basal area dervied at field survey.
- ismallg23 indicator for small area 2 and 3 for each observation.
- ismallold indicator for small area 1 for each observation.

Source

Data provided by D.Mandallaz

References

- Mandallaz, Daniel (1991). *A unified approach to sampling theory for forest inventory based on infinite population and superpopulation models*. <http://dx.doi.org/10.3929/ethz-a-000585900>
- Mandallaz, Daniel (1993). *Geostatistical methods for double sampling schemes. application to combined forest inventories*. Chair of Forest Inventory and Planning, Swiss Federal Institute of Technology (ETH). <http://dx.doi.org/10.3929/ethz-a-000943897>

Index

*Topic **datasets**

grisons, [6](#)

zberg, [23](#)

colnames, [11](#), [18](#)

confint, [2](#)

data.frame, [12](#), [19](#)

forestinventory, [4](#)

forestinventory-package
(forestinventory), [4](#)

formula, [7](#), [10](#), [17](#)

grisons, [6](#)

list, [7](#), [8](#), [10](#), [17](#)

lm, [10](#), [17](#)

logical, [10](#), [11](#), [18](#)

onephase, [3](#), [5](#), [7](#)

row.names, [11](#), [18](#)

summary, [9](#)

threephase, [3](#), [5–7](#), [10](#), [10](#), [24](#)

twophase, [3](#), [5](#), [6](#), [10](#), [17](#)

zberg, [23](#)