The R 'sampling' package

Alina Matei and Yves Tillé University of Neuchâtel

> Cardiff, Q2006 April 2006

æ

Alina Matei and Yves Tillé University of Neuchâtel The R 'sampling' package

The R language History of the package

(日) (同) (三) (三)

臣

The R language

Shareware available on http://cran.r-project.org/

Alina Matei and Yves Tillé University of Neuchâtel The R 'sampling' package

The R language History of the package

æ

- Shareware available on http://cran.r-project.org/
- The Comprehensive R Archive Network

The R language History of the package

・ロト ・ 日 ・ ・ ヨ ・ ・ ヨ ・

æ

- Shareware available on http://cran.r-project.org/
- The Comprehensive R Archive Network
- Installation: 10 minutes

The R language History of the package

(D) (A) (A)

- Shareware available on http://cran.r-project.org/
- The Comprehensive R Archive Network
- Installation: 10 minutes
- All the manuels are available in pdf

The R language History of the package

イロト イヨト イヨト イヨト

- Shareware available on http://cran.r-project.org/
- The Comprehensive R Archive Network
- Installation: 10 minutes
- All the manuels are available in pdf
- Everyone can write an additional package (600 packages are available)

The R language History of the package

イロト イポト イヨト イヨト

- Shareware available on http://cran.r-project.org/
- The Comprehensive R Archive Network
- Installation: 10 minutes
- All the manuels are available in pdf
- Everyone can write an additional package (600 packages are available)
- Packages are loaded directly from R.

The R language History of the package

イロト イポト イヨト イヨト

- Shareware available on http://cran.r-project.org/
- The Comprehensive R Archive Network
- Installation: 10 minutes
- All the manuels are available in pdf
- Everyone can write an additional package (600 packages are available)
- Packages are loaded directly from R.
- ► The manual of the package is available online and in pdf.

The R language History of the package

・ロン ・四と ・ヨン ・ヨン

- Shareware available on http://cran.r-project.org/
- The Comprehensive R Archive Network
- Installation: 10 minutes
- All the manuels are available in pdf
- Everyone can write an additional package (600 packages are available)
- Packages are loaded directly from R.
- The manual of the package is available online and in pdf.
- Package 'sampling' written by Matei and Tillé.

The R language History of the package

(日) (同) (三) (三)

æ

Continuous distributions

► EFTA course for public statisticians (April 2005).

Alina Matei and Yves Tillé University of Neuchâtel The R 'sampling' package

The R language History of the package

(D) (A) (A) (A)

- ► EFTA course for public statisticians (April 2005).
- Objective : to apply directly the theory with the R language.

The R language History of the package

- ► EFTA course for public statisticians (April 2005).
- Objective : to apply directly the theory with the R language.
- Theory + Exercices with a laptop and R.

The R language History of the package

- ► EFTA course for public statisticians (April 2005).
- Objective : to apply directly the theory with the R language.
- Theory + Exercices with a laptop and R.
- Writing of a large set of procedures.

The R language History of the package

- ► EFTA course for public statisticians (April 2005).
- ► Objective : to apply directly the theory with the R language.
- Theory + Exercices with a laptop and R.
- Writing of a large set of procedures.
- ► Finally, decision of submitting the package to the CRAN.

Topics Tools Data bases Simple random sampling Unequal probability sampling Balanced sampling

・ロト ・ 同ト ・ ヨト ・ ヨト

臣

Content of the package

Stratification, two-stage, unequal probabilities, balanced sampling

Topics Tools Data bases Simple random sampling Unequal probability sampling Balanced sampling

・ロト ・ 同ト ・ ヨト ・ ヨト

э

Content of the package

- Stratification, two-stage, unequal probabilities, balanced sampling
- Estimation: calibration and regression estimator

Topics Tools Data bases Simple random sampling Unequal probability sampling Balanced sampling

イロト イヨト イヨト イヨト

Content of the package

- Stratification, two-stage, unequal probabilities, balanced sampling
- Estimation: calibration and regression estimator
- ► Tools : computation of inclusion probabilities, crossing strata

Topics Tools Data bases Simple random sampling Unequal probability sampling Balanced sampling

Content of the package

- Stratification, two-stage, unequal probabilities, balanced sampling
- Estimation: calibration and regression estimator
- ► Tools : computation of inclusion probabilities, crossing strata
- Data bases, Swiss municipalities, Belgian municipalities.

Topics Tools Data bases Simple random sampling Unequal probability sampling Balanced sampling

・ロト ・ 日 ・ ・ ヨ ・ ・ ヨ ・

臣



writesample: return the list of all the samples of fixed sample size

Alina Matei and Yves Tillé University of Neuchâtel The R 'sampling' package

Topics Tools Data bases Simple random sampling Unequal probability sampling Balanced sampling

・ロト ・日本 ・モート ・モート



- writesample: return the list of all the samples of fixed sample size
- cleanstrata: renumbering of the strata

Alina Matei and Yves Tillé University of Neuchâtel The R 'sampling' package

Topics Tools Data bases Simple random sampling Unequal probability sampling Balanced sampling

ヘロン ヘロン ヘヨン ヘ



- writesample: return the list of all the samples of fixed sample size
- cleanstrata: renumbering of the strata
- disjonctive return a matrix wit 0 and 1 that is the disjonctive representation of the stratum.

Topics Tools Data bases Simple random sampling Unequal probability sampling Balanced sampling

ヘロン ヘロン ヘヨン ヘ



- writesample: return the list of all the samples of fixed sample size
- cleanstrata: renumbering of the strata
- disjonctive return a matrix wit 0 and 1 that is the disjonctive representation of the stratum.
- inclusionprobabilities: compute unequal inclusion probabilities from an auxiliary variable variable.

Tools Data bases Simple random sampling Unequal probability sampling Balanced sampling

・ロト ・日本 ・モート ・モート



 MU284 A data frame with 284 municipalities on the following 11 variables : populations, political results.

Topics Tools **Data bases** Simple random sampling Unequal probability sampling Balanced sampling

(D) (A) (A) (A)

Data bases

- MU284 A data frame with 284 municipalities on the following 11 variables : populations, political results.
- swissmunicipalities: 2896 Swiss municipalities. Surfaces and population.

Topics Tools **Data bases** Simple random sampling Unequal probability sampling Balanced sampling

Data bases

- MU284 A data frame with 284 municipalities on the following 11 variables : populations, political results.
- swissmunicipalities: 2896 Swiss municipalities. Surfaces and population.
- belgianmunicipalities: 589 Belgian municipalities 11 variables, population and taxes.

Topics Tools Data bases Simple random sampling Unequal probability sampling Balanced sampling

・ロト ・ 同ト ・ ヨト ・ ヨト

臣

Simple random sampling

srswor: Simple random sampling with replacement.

Alina Matei and Yves Tillé University of Neuchâtel The R 'sampling' package

Topics Tools Data bases Simple random sampling Unequal probability sampling Balanced sampling

イロト イポト イヨト イヨト

Simple random sampling

- srswor: Simple random sampling with replacement.
- srswor1: Simple random sampling without replacement (sequential method).

Topics Tools Data bases Simple random sampling Unequal probability sampling Balanced sampling

Simple random sampling

- srswor: Simple random sampling with replacement.
- srswor1: Simple random sampling without replacement (sequential method).
- srswr: Simple random sampling with replacement.

Topics Tools Data bases Simple random sampling **Unequal probability sampling** Balanced sampling

・ロト ・日本 ・モート ・モート

臣

Unequal probability sampling

UPbrewer,

Alina Matei and Yves Tillé University of Neuchâtel The R 'sampling' package

Topics Tools Data bases Simple random sampling Unequal probability sampling Balanced sampling

æ

Unequal probability sampling

UPbrewer,

UPmaxentropy, (set of function)

Topics Tools Data bases Simple random sampling Unequal probability sampling Balanced sampling

・ロト ・日下・ ・日下

- UPbrewer,
- UPmaxentropy, (set of function)
- UPmidzuno, UPmidzunopi2,

Topics Tools Data bases Simple random sampling Unequal probability sampling Balanced sampling

A = A = A
 A = A
 A = A
 A
 A = A
 A
 A = A
 A
 A = A
 A
 A = A
 A
 A = A
 A
 A = A
 A
 A = A
 A
 A = A
 A
 A = A
 A
 A = A
 A
 A = A
 A
 A = A
 A
 A = A
 A
 A = A
 A
 A = A
 A
 A = A
 A
 A = A
 A
 A = A
 A
 A = A
 A
 A = A
 A
 A = A
 A
 A = A
 A
 A = A
 A
 A = A
 A
 A = A
 A
 A = A
 A
 A = A
 A
 A = A
 A
 A = A
 A
 A
 A = A
 A
 A
 A = A
 A
 A
 A = A
 A
 A
 A
 A
 A = A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A

- UPbrewer,
- UPmaxentropy, (set of function)
- UPmidzuno, UPmidzunopi2,
- UPmultinomial,

Topics Tools Data bases Simple random sampling Unequal probability sampling Balanced sampling

- UPbrewer,
- UPmaxentropy, (set of function)
- UPmidzuno, UPmidzunopi2,
- UPmultinomial,
- UPpivotal, UPrandompivotal,

Topics Tools Data bases Simple random sampling Unequal probability sampling Balanced sampling

(D) (A) (A) (A)

- UPbrewer,
- UPmaxentropy, (set of function)
- UPmidzuno, UPmidzunopi2,
- UPmultinomial,
- UPpivotal, UPrandompivotal,
- UPpoisson,

Topics Tools Data bases Simple random sampling Unequal probability sampling Balanced sampling

- E - N

- UPbrewer,
- UPmaxentropy, (set of function)
- UPmidzuno, UPmidzunopi2,
- UPmultinomial,
- UPpivotal, UPrandompivotal,
- UPpoisson,
- UPsampford,

Topics Tools Data bases Simple random sampling **Unequal probability sampling** Balanced sampling

・ロト ・ 日 ・ ・ ヨ ・ ・ ヨ ・

- UPbrewer,
- UPmaxentropy, (set of function)
- UPmidzuno, UPmidzunopi2,
- UPmultinomial,
- UPpivotal, UPrandompivotal,
- UPpoisson,
- UPsampford,
- UPsystematic, UPrandomsystematic, UPsystematicpi2,

Topics Tools Data bases Simple random sampling **Unequal probability sampling** Balanced sampling

- UPbrewer,
- UPmaxentropy, (set of function)
- UPmidzuno, UPmidzunopi2,
- UPmultinomial,
- UPpivotal, UPrandompivotal,
- UPpoisson,
- UPsampford,
- UPsystematic, UPrandomsystematic, UPsystematicpi2,
- ▶ UPtille, UPtillepi2,

Topics Tools Data bases Simple random sampling Unequal probability sampling Balanced sampling

・ロト ・ 同ト ・ ヨト ・ ヨト

æ

Balanced sampling

Design that satisfies the balancing equations

$$\sum_{k\in S}\frac{\mathbf{x}_k}{\pi_k}=\sum_{k\in U}\mathbf{x}_k,$$

where \mathbf{x}_k is a vector of auxiliary variables.

Topics Tools Data bases Simple random sampling Unequal probability sampling Balanced sampling

イロト イポト イヨト イヨト

Balanced sampling

Design that satisfies the balancing equations

$$\sum_{k\in S}\frac{\mathbf{x}_k}{\pi_k}=\sum_{k\in U}\mathbf{x}_k,$$

where \mathbf{x}_k is a vector of auxiliary variables.

Cube algorithm: flight phase and landing phase.

Topics Tools Data bases Simple random sampling Unequal probability sampling Balanced sampling

Balanced sampling

Design that satisfies the balancing equations

$$\sum_{k\in S}\frac{\mathbf{x}_k}{\pi_k}=\sum_{k\in U}\mathbf{x}_k,$$

where \mathbf{x}_k is a vector of auxiliary variables.

- Cube algorithm: flight phase and landing phase.
- samplecube, fastflightcube, landingcube

Topics Tools Data bases Simple random sampling Unequal probability sampling Balanced sampling

・ロト ・ 日 ・ ・ ヨ ・ ・ ヨ ・

Balanced sampling

Design that satisfies the balancing equations

$$\sum_{k\in S}\frac{\mathbf{x}_k}{\pi_k}=\sum_{k\in U}\mathbf{x}_k,$$

where \mathbf{x}_k is a vector of auxiliary variables.

- Cube algorithm: flight phase and landing phase.
- samplecube, fastflightcube, landingcube
- Complex survey balancedstratification balancedcluster balancedtwostage

Topics Tools Data bases Simple random sampling Unequal probability sampling Balanced sampling

ヘロン ヘロン ヘヨン ヘ



Exercise

Compute inclusion probabilities 200 Belgian municipalities with inclusion probabilities proportional to the population in 2004.

Topics Tools Data bases Simple random sampling Unequal probability sampling Balanced sampling

Exercises

Exercise

Use the Belgian database. Select a sample of 200 municipalities with unequal probabilities proportional to the number of inhabitants in 2004.

- with Poisson sampling
- with a method of unequal probabilities and fixed sample size
- with simple random sampling.

Compute the Horvitz-Thompson estimators of the taxable income for 50 samples and draw a boxplot of the estimators for each method.

Topics Tools Data bases Simple random sampling Unequal probability sampling Balanced sampling

・ロン ・四と ・ヨン ・ヨン

Exercises

Exercise

Use the database of Swiss municipalities, and select a stratified balanced sample. A balanced sample is first selected in each strata. Next the results of the flight phase are merged and a flight phase is applied again on the whole population. Finally, a landing phase is applied on all the population. Use the following balancing variables: HApoly, Surfacesbois, POOBMTOT, POOBWTOT, POPTOT, Pop020, Pop2040, Pop4065, Pop65P, H00PTOT. The sample size is 400 and the municipalities must be selected with inclusion probabilities proportional to POPTOT. The stratification variable is REG (swiss regions). Next, print the names of the selected municipalities.

Topics Tools Data bases Simple random sampling Unequal probability sampling Balanced sampling

Exercises

Exercise

Use the Belgian database. Select a sample of 200 municipalities with unequal probabilities proportional to the number of inhabitants in 2004 with Poisson sampling design. Next calibrate the sample by means of the raking ratio estimator on the variables: mean(Men03), mean(Women03), Diffmen, Diffwom, TaxableIncome, Totaltaxation, averageincome, medianincome. The division by the means is necessary to avoid too large numbers. Compute the Horvitz-Thompson estimators and the calibrated estimators for the calibration variables. Limit the variation of the g-weights between 0.5 and 1.5.