

Survey Sampling Methodology in Latvia

Mārtiņš Liberts

Central Statistical Bureau of Latvia

8-12 December 2013

Mārtiņš Liberts (CSB)



Introduction

Software

Sampling Design

Non-response and Weighting

Imputation

Sampling Errors

Organisation of Methodological Work

Mārtiņš Liberts (CSB)



Introduction

Software

Sampling Design

Non-response and Weighting

Imputation

Sampling Errors

Organisation of Methodological Work



- ► Established on the 1st September 1919
- ► Incorporation in statistical system of Soviet Union 1945
- ► Independence regained in 1991
- ► 549 employees at the beginning of 2013
- ► The main provider of the official statistics in Latvia
- ► Survey methodology is used since 90-ties.

Mathematical Support Division

- Survey methodology is applied in centralised manner
- ► Mathematical Support Division is responsible for:
 - Sampling design and sampling
 - Weighting of survey data
 - Imputation (only for social surveys)
 - Precision estimation (sampling errors)



Mathematical Support Division

- Survey methodology is applied in centralised manner
- ► Mathematical Support Division is responsible for:
 - Sampling design and sampling
 - Weighting of survey data
 - Imputation (only for social surveys)
 - Precision estimation (sampling errors)
 - ► Time series analyses:
 - Seasonal adjustment
 - Short term forecasting
 - Training of the CSB staff





Introduction

Software

Sampling Design

Non-response and Weighting

Imputation

Sampling Errors

Organisation of Methodological Work

Software Used for Survey Methodology



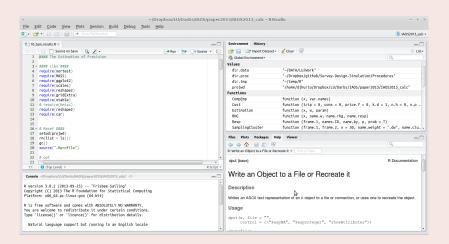
- ► SPSS
 - Sampling
 - Weighting
 - Imputation
- R (http://www.r-project.org/)
 RStudio (http://www.rstudio.com/)
 - Sampling
 - ► Weighting, calibration
 - Imputations
 - Sampling error estimation
- ► Demetra+, JDemetra+
 - Seasonal Adjustment



- ► Sampling
 - function dom_optimal_allocation
 - ► Jānis Jukāms, Central Statistical Bureau of Latvia
- Calibrations
 - package sampling function calib
 - Yves Tillé and Alina Matei (2012). sampling: Survey Sampling. R package version 2.5. http://CRAN.R-project.org/package=sampling
- Sampling error estimation
 - package vardpoor function vardom
 - Juris Breidaks, Mārtiņš Liberts (2013) Central Statistical Bureau of Latvia



- Sampling error estimation
 - package vardpoor function vardom
 - Juris Breidaks, Mārtiņš Liberts (2013) Central Statistical Bureau of Latvia



~/Dropbox/LU/Darbs/IAOS/paper2013/IAOS2013 calc - RStudio File Edit Code View Plots Session Build Debug Tools Help 🔍 🔹 🗲 📲 🔚 🔚 🔚 🚺 🏕 Co to file/function IAOS2013 calc * 10_Sym_results.R × Environment History 📄 🔲 Source on Save 🛛 🔍 🧪 🗸 A Run SA Source * 😭 📊 🔄 Import Dataset + 🥑 Clear 🕓 ≡ List+ 1 #### The Estimation of Precision Global Environment -3 - #### Libs #### About RStudio 4 require(nortest) 5 require(MASS) lthub/Survey-Design-Simulation/Procedures" RStudio 6 require(opplot2) 7 require(scales) Version 0.98.484 - © 2009-2013 RStudio. Inc. lo/Dropbox/LU/Darbs/IAOS/paper2013/IAOS2013_calc" 8 require(reshape2) 9 require(gridExtra) Mozilla/5.0 (X11: Linux x86 64) AppleWebKit/534.34 (KHTML, like Gecko) RStudio Safari/534.34 Qt/4.8.0 10 require(xtable) var.names) 11 # require(Hmisc) Unless you have received this program directly from RStudio pursuant to the terms of a ip = 0, cons = 0, price.f = 0, k.d = 1, n.h = 0, n.p .. 12 require(reshape2) commercial license agreement with RStudio, then this program is licensed to you under the terms w, paran) of version 3 of the GNU Affero General Public License. 13 require(car) name.w, name.rhg, name.resp) 14 ame.1, names.ID, name.by, p, prob = T) RStudio includes other open source software components. The following is a list of these components (full copies of the license agreements ame.1, frame.2, n = 30, name.weight = ".dw", name.clu. 16 - # Reset #### used by these components are included below): 17 setwd(projwd) 18 rn(list = ls()) Qt (LGPL v2.1) OtSingleApplication 19 ac() Ace (LGPL v2.1) 20 source(".Rprofile") 21 RapidXnl 22 # opt JSON Spirit R Documentation Google Web Toolkit 🚺 (Top Level) 🌣 GIN or Recreate it AOP Alliance Console -/Dropbox/LU/Darbs/IAOS/paper2013/IAOS2013_calc/ 🔗 RSA-35 R version 3.0.2 (2013-09-25) -- "Frisbee Sailing" Copyright (C) 2013 The R Foundation for Statistical Com OK Platform: x86 64-pc-linux-gnu (64-bit) ext to a file or connection, or uses one to recreate the object. R is free software and comes with ABSOLUTELY NO WARRANTY. Usage You are welcome to redistribute it under certain conditions. Type 'license()' or 'licence()' for distribution details. dput(x, file = "", control = c("keepNA", "keepInteger", "showAttributes")) Natural language support but running in an English locale ----



Introduction

Software

Sampling Design

Non-response and Weighting

Imputation

Sampling Errors

Organisation of Methodological Work

Mārtiņš Liberts (CSB)

Sampling Frame in Household Surveys



- ► Source: Statistical Dwelling Register
 - Population Register
 - Building Register
 - ► Address Register
 - ► other
- Output for sampling:
 - List of census counting areas
 - List of dwellings
 - ► List of persons

- Sampling design depends on mode of survey:
 - ► Paper assisted **personal** interviews (PAPI)
 - ► Computer assisted **personal** interviews (CAPI)
 - ► Computer assisted telephone interviews (CATI)
 - ► Computer assisted web interviews (CAWI)
- CAPI is a traditional mode
- ► The usage of CATI is increasing
- ► CAPI or CAPI/CATI are the most common modes
- CAWI was used in the last population census (2011)



- ► There are travelling costs if CAPI is used
- ► We want to minimize / optimise the travelling costs
- Two-stage sampling design is used:
 - ► Sampling of (*census counting*) areas is used in the 1st stage
 - ► Sampling of dwellings or individuals is used in the 2nd stage
- ► This allows to reduce / control travelling costs



- ► There are travelling costs if CAPI is used
- ► We want to minimize / optimise the travelling costs
- Two-stage sampling design is used:
 - ► Sampling of (*census counting*) areas is used in the 1st stage
 - \blacktriangleright Sampling of dwellings or individuals is used in the 2nd stage
- ► This allows to reduce / control travelling costs
- Attention: The sampling errors tend to increase because of a clustering effect

- Stratified simple random sampling (SSRS)
 - Survey on doctoral degree holders (very small population)
- Two-stage sampling for CAPI and SSRS for CATI
 - European Health Survey (2014)



- ► Source: Statistical Business Register
 - ► State Business Register
 - State Revenue Service (tax office)
 - ► other
- Output for sampling:
 - List of active enterprises and organisations



- ► CAWI (e-questionnaire) is the most common mode
- postal surveys
- call back is used in all cases if necessary



- ► CAWI (e-questionnaire) is the most common mode
- postal surveys
- call back is used in all cases if necessary
- Remark: There are no travelling costs



- ► Stratified simple random sampling is used in most cases
- Stratification variables:
 - ► Size groups (size measured by turnover, number of employees)
 - Economic activity branch (NACE classification)
 - ► Type of unit
 - Region
- ► Optimal sample allocation for each domain (*R procedure*)



- Two stage sampling for the Survey on Employees:
 - ► Units (enterprises or local units) sampled in the first stage
 - Employees sampled in the second stage

Sampling Frame in Agriculture Surveys



- ► Source: Statistical Farm Register
 - ► State Land Service
 - ► Farming Land Register
 - ► Animal Register
 - ► other
- Output for sampling:
 - List of active farms



- ► CAWI (e-questionnaire) is used as the first mode
- CAPI is used as the second mode



- ► CAWI (e-questionnaire) is used as the first mode
- CAPI is used as the second mode
- ► Remark: There are travelling costs (in CAPI mode)



- ► Stratified simple random sampling is used in most cases
- Stratification variables:
 - ► Size groups (size measured by land area, economic size)
 - Specialisation
 - Region



Introduction

Software

Sampling Design

Non-response and Weighting

Imputation

Sampling Errors

Organisation of Methodological Work

and the second

Usual scheme:

- Design weights according to sampling design (computed during sampling)
- ► Non-response correction by response homogeneity groups
- Trimming of extreme weights (optional)
- Calibration of weights to external information (population counts)
 - ▶ Function calib() from the sampling package (R)



 Cross-sectional and longitudinal weighting for European Survey on Income and Living Conditions (EU-SILC)



Usual scheme:

- Design weights according to sampling design (computed during sampling)
- Usually population frame is updated (sampling frame and weighting frame)
- Weighting to population counts:
 - Number raised estimator in each stratum
 - Calibration of weights



Introduction

Software

Sampling Design

Non-response and Weighting

Imputation

Sampling Errors

Organisation of Methodological Work

- Usually imputation is done for income variables
- Methods:
 - Randomised hot-deck imputation in groups
 - ► Nearest neighbour (distance function) imputation
 - Regression:
 - Linear regression
 - Interval regression



- Done by subject matter unit (decentralised)
- Case specific (large enterprises are heterogeneous)
- Imputation using:
 - Administrative records
 - Historical data



Introduction

Software

Sampling Design

Non-response and Weighting

Imputation

Sampling Errors

Organisation of Methodological Work

Mārtiņš Liberts (CSB)

Estimation of Sampling Error

and the second

- Common procedure for all surveys
- Written as an R procedure function vardom from package vardpoor
- Main steps of the procedure:
 - ► Domain variables are generated for domain estimates

$$y_{i,d} = \begin{cases} y_i & \text{if } i \in D_d \\ 0 & \text{if } i \notin D_d \end{cases}$$

- Linearised variables are computed for non-linear statistics (ratio of two totals, Gini index, ...)
- Residuals from the regression model are estimated if weight calibration is applied
- Ultimate cluster variance estimator (Hansen, Hurwitz and Madow, 1953)



Ultimate cluster variance estimator

- Ultimate cluster variance estimator (Hansen, Hurwitz and Madow, 1953)
- Osier (2012) "The Linearisation Approach Implemented by Eurostat for The First Wave of EU-SILC: What Could Be Done from The Second Wave Onwards?"

$$\hat{V}\left(\hat{\Theta}\right) = \sum_{h=1}^{H} \left(1 - \frac{n_h}{N_h}\right) \frac{n_h}{n_h - 1} \sum_{i=1}^{n_h} \left(z_{hi\bullet} - \bar{z}_{h\bullet\bullet}\right)^2$$

$$z_{hi\bullet} = \sum_{j=1}^{m_{hi}} \omega_{hij} \cdot z_{hij}$$
$$\bar{z}_{h\bullet\bullet} = \frac{\sum_{i=1}^{n_h} z_{hi\bullet}}{n_h}$$

Input for vardom



- Y study variables
- H stratification
- PSU ID of primary sampling units
- w_final final weight
- Dom domain variables
- N_h PSU population size in each stratum
- Z denominator variables (for ratio)
- ► X calibration variables
- ▶ g calibration factor (g-weight)



- ▶ estim parameter estimator
- var variance
- se standard error
- ▶ cv coefficient of variation
- CI_lower the lower bound if confidence interval
- CI_upper the upper bound if confidence interval
- deff design effect



Introduction

Software

Sampling Design

Non-response and Weighting

Imputation

Sampling Errors

Organisation of Methodological Work



- Information Technology Division updates the Statistical Dwelling Register (SDR)
- Mathematical Support Division (MSD) extracts information from SDR to build the sampling frame (SPSS syntax)
- MSD creates sample file
- Sample file is sent to:
 - Interview Organisation Division
 - Employment Statistics Division (Labour Force Survey)
 - Income Statistics Division (Survey on Income and Living Standards)



- Business Register Division (BRD) updates the Statistical Business Register
- Mathematical Support Division (MSD) receives a file from BRD with active units (enterprises, organisations)
- MSD creates sample file
- Sample file is sent to subject matter divisions

- Information Technology Division (ITD) updates the Statistical Farm Register
- Mathematical Support Division (MSD) receives a file from ITD with active farms
- MSD creates sample file
- ► Sample file is sent to Agriculture Statistics Division



Thank you!

Mārtiņš Liberts (CSB)