

Package ‘StatMatch’

January 8, 2025

Version 1.4.3

Date 2025-01-08

Title Statistical Matching or Data Fusion

Author Marcello D'Orazio [aut, cre]

Maintainer Marcello D'Orazio <mdo.statmatch@gmail.com>

Depends R (>= 2.7.0), proxy, survey, lpSolve, ggplot2, dplyr

Suggests Hmisc, MASS, mipfp, clue, RANN, randomForest, nnet, glmnet, naivebayes

Description Integration of two data sources referred to the same target population which share a number of variables. Some functions can also be used to impute missing values in data sets through hot deck imputation methods. Methods to perform statistical matching when dealing with data from complex sample surveys are available too.

License GPL (>= 2)

URL <https://github.com/marcellodo/StatMatch>,
https://github.com/marcellodo/StatMatch/tree/master/Tutorials_Vignette_OtherDocs

NeedsCompilation no

Repository CRAN

Date/Publication 2025-01-08 22:00:02 UTC

Contents

| | |
|-----------------------------|----|
| StatMatch-package | 2 |
| comb.samples | 3 |
| comp.conf | 8 |
| comp.prop | 10 |
| create.fused | 13 |
| create.imputed | 15 |
| fact2dummy | 17 |
| Fbounds.pred | 18 |
| Fbwidths.by.x | 21 |

| | |
|------------------------------|----|
| Frechet.bounds.cat | 24 |
| gower.dist | 28 |
| harmonize.x | 31 |
| mahalanobis.dist | 36 |
| maximum.dist | 37 |
| mixed.mtc | 39 |
| NND.hotdeck | 44 |
| pBayes | 48 |
| plotBounds | 51 |
| plotCont | 52 |
| plotTab | 54 |
| pw.assoc | 55 |
| RANDwNND.hotdeck | 58 |
| rankNND.hotdeck | 63 |
| rho.bounds | 66 |
| rho.bounds.pred | 68 |
| samp.A | 71 |
| samp.B | 72 |
| samp.C | 73 |
| selMtc.by.unc | 75 |

| | |
|--------------|-----------|
| Index | 79 |
|--------------|-----------|

| | |
|-------------------|--|
| StatMatch-package | <i>Statistical Matching or Data Fusion</i> |
|-------------------|--|

Description

Functions to perform statistical matching (aka data fusion), i.e. the integration of two data sources. Some functions can also be used to impute missing values in data sets through hotdeck imputation methods.

Details

Statistical matching (also known as *data fusion*) aims to integrate two data sources (typically samples) that refer to the same target population and share a number of variables. The ultimate aim is to study the relationship between variables that are not jointly observed in a single data source. The integration can be performed at *micro* level (a *synthetic* or *fused* file is the output) or *macro* level (the aim is to estimate correlation coefficients, regression coefficients, contingency tables, etc.). Non-parametric hotdeck imputation methods (random, rank and nearest neighbour) can be used to derive the synthetic data set. Alternatively, a mixed (parametric-nonparametric) procedure based on predictive mean matching can be used. Methods are also available for performing statistical matching when dealing with data from complex sample surveys. Finally, some functions can be used to explore the uncertainty due to the typical matching framework and can help to decide whether or not to proceed with statistical matching at the micro or macro level. For major details see D’Orazio et al. (2006), D’Orazio (2015), and material available at https://github.com/marcellodo/StatMatch/tree/master/Tutorials_Vignette_OtherDocs.

Author(s)

Marcello D'Orazio

Maintainer: Marcello D'Orazio <mdo.statmatch@gmail.com>

References

D'Orazio M. (2015) Integration and imputation of survey data in R: the StatMatch package, *Romanian Statistical Review*, 2/2015, pp. 57–68

D'Orazio M., Di Zio M., Scanu M. (2006) *Statistical Matching, Theory and Practice*. Wiley, Chichester.

comb.samples

Statistical Matching of data from complex sample surveys

Description

This function allows you to cross-tabulate two categorical variables, Y and Z, observed separately in two independent surveys (Y is collected in survey A and Z is collected in survey B) relating to the same target population. The two surveys share a number of common variables **X**. If a third survey C is available on the same population and collects both Y and Z, these data are used as a source of additional information.

The statistical adjustment is done by calibrating the survey weights as suggested in Renssen (1998).

It is also possible to use the function to derive estimates that a unit falls into one of the categories of the target variable (the estimates are based on Linear probability models and are obtained as a by-product of the Renssen method).

Usage

```
comb.samples(svy.A, svy.B, svy.C=NULL, y.lab, z.lab, form.x,  
             estimation=NULL, micro=FALSE, ...)
```

Arguments

- | | |
|-------|--|
| svy.A | A svydesign R object that stores the data collected in the survey A and all the information concerning the corresponding sampling design. This object can be created by using the function svydesign in the package survey . All the variables specified in form.x and by y.lab must be available in svy.A. |
| svy.B | A svydesign R object that stores the data collected in the survey B and all the information concerning the corresponding sampling design. This object can be created by using the function svydesign in the package survey . All the variables specified in form.x and by z.lab must be available in svy.B. |
| svy.C | A svydesign R object that stores the data collected in the the survey C and all the information concerning the corresponding sampling design. This object can be created by using the function svydesign in the package survey . |

| | |
|-------------------------|--|
| | <p>When <code>svy.C=NULL</code> (default), i.e. no auxiliary information is available, the function returns an estimate of the contingency table of <i>Y</i> vs. <i>Z</i> under the Conditional Independence assumption (CIA) (see Details for major information).</p> <p>When <code>svy.C</code> is available, if <code>estimation="incomplete"</code> then it must contain at least <code>y.lab</code> and <code>z.lab</code> variables. On the contrary, when <code>estimation="synthetic"</code> all the variables specified in <code>form.x</code>, <code>y.lab</code> and <code>z.lab</code> must be available in <code>svy.C</code>.</p> |
| <code>y.lab</code> | A string providing the name of the <i>Y</i> variable, available in survey A and in survey C (if available). The <i>Y</i> variable can be a categorical variable (factor in R) or a continuous one (in this latter case <code>z.lab</code> should be categorical). |
| <code>z.lab</code> | A string providing the name of the <i>Z</i> variable available in survey B and in survey C (if available). The <i>Z</i> variable can be a categorical variable (factor in R) or a continuous one (in this latter case <code>y.lab</code> should be categorical). |
| <code>form.x</code> | <p>A R formula specifying which matching variables (subset of the X variables) are collected in all surveys and how they are to be considered when combining samples. For example, <code>form.x=~x1+x2</code> means that variables <code>x1</code> and <code>x2</code> are to be considered marginally, without taking into account their cross-tabulation; only their marginal distribution is considered. To skip the intercept, the formula must be written as <code>form.x=~x1+x2-1</code>.</p> <p>When dealing with categorical variables, <code>form.x=~x1:x2-1</code> means that the joint distribution of the two variables (table of <code>x1</code> vs. <code>x2</code>) has to be considered.</p> <p>To better understand the use of <code>form.x</code>, see model.matrix (see also formula).</p> <p>Due to the weight calibration features, it is preferable to work with categorical X variables. In some cases, the procedure can be successful when a single continuous variable is considered together with one or more categorical variables, but often it may be necessary to categorise the continuous variable (see Details).</p> |
| <code>estimation</code> | <p>A string identifying the method to be used to estimate the table of <i>Y</i> vs. <i>Z</i> when data from survey C are available. As suggested in Renssen (1998), there are two alternative methods: (i) incomplete two-way stratification (<code>estimation="incomplete"</code> or <code>estimation="ITWS"</code>, the default) and (ii) synthetic two-way stratification (<code>estimation="synthetic"</code> or <code>estimation="STWS"</code>). In the first case (<code>estimation="incomplete"</code>), only <i>Y</i> and <i>Z</i> variables need to be available in <code>svy.C</code>. On the other hand, in the case of <code>estimation="synthetic"</code>, the C survey must contain all the X variables specified by <code>form.x</code>, as well as the <i>Y</i> and <i>Z</i> variables. See details for more information.</p> |
| <code>micro</code> | Logical, when TRUE predictions of <i>Z</i> in A and of <i>Y</i> in B are provided. In particular when <i>Y</i> and <i>Z</i> are both categorical variables it is provided the estimated probability that a unit falls in each of the categories of the given variable. These probabilities are estimated as a by-product of the whole procedure by considering Linear Probability Models, as suggested in Renssen (1998) (see Details) |
| <code>...</code> | <p>Further arguments that may be necessary for calibration. In particular, the argument <code>calfun</code> allows to specify the calibration function:</p> <ul style="list-style-type: none"> (i) <code>calfun="linear"</code> for linear calibration (default); (ii) <code>calfun="raking"</code> to rake the survey weights; and (iii) <code>calfun="logit"</code> for logit calibration. See calibrate for major details. |

Note that when `calfun="linear"` calibration may return negative weights. Generally speaking, in sample surveys weights are expected to be greater than or equal to 1, i.e. `bounds=c(1, Inf)`.

The number of iterations used in calibration can be modified by using the argument `maxit` (by default `maxit=50`).

See [calibrate](#) for further details.

Details

This function estimates the Y vs. Z contingency table by performing a series of calibrations of the survey weights. In practice, the estimation is carried out on data from survey C, using all the information from surveys A and B. If survey C is not available, the table of Y vs. Z is estimated under the assumption of conditional independence (CIA), i.e. $p(Y, Z) = p(Y|\mathbf{X}) \times p(Z|\mathbf{X}) \times p(\mathbf{X})$.

When data from survey C are available (Renssen, 1998), the table of Y vs. Z can be estimated by: Incomplete Two-Way Stratification (ITWS) or Synthetic Two-Way Stratification (STWS). In the first case (ITWS) the weights of the units in survey C are calibrated so that the new weights allow to reproduce the marginal distributions of Y estimated on survey A, and that of Z estimated on survey B. Note that the distribution of the X variables in survey A and in survey B, must be harmonized before performing ITWS (see [harmonize.x](#)).

The Synthetic Two-Way Stratification allows to estimate the table of Y vs. Z by considering also the X variables observed in C. This method consists in correcting the table of Y vs. Z estimated under the CIA according to the relationship between Y and Z observed in survey C (for major details see Renssen, 1998).

If the argument `micro` is set to TRUE the function will also return Z.A and Y.B. The first data.frame has the same rows as `svy.A` and the number of columns is equal to the number of categories of the Z variable specified by `z.lab`. Each row gives the estimated probabilities of taking a value in each category. The same happens for Y.B, which gives the estimated probabilities of assuming a category of `y.lab` for each unit in B. The estimated probabilities are obtained by applying the linear probability models (for more details see Renssen, 1998). Unfortunately, such models may give estimated probabilities less than 0 or greater than 1. Great care should be taken when using such predictions for practical purposes.

Value

A R list with the results of the calibration procedure according to the input arguments.

| | |
|---------------------|---|
| <code>yz.CIA</code> | The table of Y (<code>y.lab</code>) vs. Z (<code>z.lab</code>) estimated under the Conditional Independence Assumption (CIA). |
| <code>cal.C</code> | The survey object <code>svy.C</code> after the calibration. Only when <code>svy.C</code> is provided. |
| <code>yz.est</code> | The table of Y (<code>y.lab</code>) vs. Z (<code>z.lab</code>) estimated under the method specified via estimation argument. Only when <code>svy.C</code> is provided. |
| <code>Z.A</code> | Only when <code>micro=TRUE</code> . It is a data frame with the same rows as in <code>svy.A</code> and the number of columns is equal to the number of categories of the variable <code>z.lab</code> . Each row provides the estimated probabilities for a unit being in the various categories of <code>z.lab</code> . |

| | |
|------|---|
| Y.B | Only when micro=TRUE. It is a data frame with the same rows as in svy.B and the number of columns is equal to the number of categories of y.lab. Each row provides the estimated probabilities for a unit being in the various categories of y.lab. |
| call | Stores the call to this function with all the values specified for the various arguments (call=match.call()). |

Author(s)

Marcello D’Orazio <mdo.statmatch@gmail.com>

References

D’Orazio, M., Di Zio, M. and Scanu, M. (2006). *Statistical Matching: Theory and Practice*. Wiley, Chichester.

Renssen, R.H. (1998) “Use of Statistical Matching Techniques in Calibration Estimation”. *Survey Methodology*, **24**, pp. 171–183.

See Also

[calibrate](#), [svydesign](#), [harmonize.x](#)

Examples

```
data(quine, package="MASS") #loads quine from MASS
str(quine)
quine$c.Days <- cut(quine$Days, c(-1, seq(0,20,10),100))
table(quine$c.Days)

# split quine in two subsets

suppressWarnings(RNGversion("3.5.0"))
set.seed(124)
lab.A <- sample(nrow(quine), 70, replace=TRUE)
quine.A <- quine[lab.A, c("Eth", "Sex", "Age", "Lrn")]
quine.B <- quine[-lab.A, c("Eth", "Sex", "Age", "c.Days")]

# create svydesign objects
require(survey)
quine.A$f <- 70/nrow(quine) # sampling fraction
quine.B$f <- (nrow(quine)-70)/nrow(quine)
svy.qA <- svydesign(~1, fpc=~f, data=quine.A)
svy.qB <- svydesign(~1, fpc=~f, data=quine.B)

# Harmonization wrt the joint distribution
# of ('Sex' x 'Age' x 'Eth')

# vector of population total known
# estimated from the full data set
# note the formula!
```

```

tot.m <- colSums(model.matrix(~Eth:Sex:Age-1, data=quine))
tot.m

out.hz <- harmonize.x(svy.A=svy.qA, svy.B=svy.qB, x.tot=tot.m,
  form.x=~Eth:Sex:Age-1, cal.method="linear")

# estimation of 'Lrn' vs. 'c.Days' under the CIA

svy.qA.h <- out.hz$cal.A
svy.qB.h <- out.hz$cal.B

out.1 <- comb.samples(svy.A=svy.qA.h, svy.B=svy.qB.h,
  svy.C=NULL, y.lab="Lrn", z.lab="c.Days",
  form.x=~Eth:Sex:Age-1)

out.1$yz.CIA
addmargins(out.1$yz.CIA)

#
# incomplete two-way stratification

# select a sample C from quine
# and define a survey object

suppressWarnings(RNGversion("3.5.0"))
set.seed(4321)
lab.C <- sample(nrow(quine), 50, replace=TRUE)
quine.C <- quine[lab.C, c("Lrn", "c.Days")]
quine.C$f <- 50/nrow(quine) # sampling fraction
svy.qC <- svydesign(~1, fpc=~f, data=quine.C)

# call comb.samples
out.2 <- comb.samples(svy.A=svy.qA.h, svy.B=svy.qB.h,
  svy.C=svy.qC, y.lab="Lrn", z.lab="c.Days",
  form.x=~Eth:Sex:Age-1, estimation="incomplete",
  calfun="linear", maxit=100)

summary(weights(out.2$cal.C))
out.2$yz.est # estimated table of 'Lrn' vs. 'c.Days'
# difference wrt the table 'Lrn' vs. 'c.Days' under CIA
addmargins(out.2$yz.est)-addmargins(out.2$yz.CIA)

# synthetic two-way stratification
# only macro estimation

quine.C <- quine[lab.C, ]
quine.C$f <- 50/nrow(quine) # sampling fraction
svy.qC <- svydesign(~1, fpc=~f, data=quine.C)

out.3 <- comb.samples(svy.A=svy.qA.h, svy.B=svy.qB.h,
  svy.C=svy.qC, y.lab="Lrn", z.lab="c.Days",
  form.x=~Eth:Sex:Age-1, estimation="synthetic",
  calfun="linear", bounds=c(.5, Inf), maxit=100)

```

```
summary(weights(out.3$cal.C))

out.3$yz.est # estimated table of 'Lrn' vs. 'c.Days'
# difference wrt the table of 'Lrn' vs. 'c.Days' under CIA
addmargins(out.3$yz.est)-addmargins(out.3$yz.CIA)
# diff wrt the table of 'Lrn' vs. 'c.Days' under incomplete 2ws
addmargins(out.3$yz.est)-addmargins(out.2$yz.CIA)

# synthetic two-way stratification
# with micro predictions

out.4 <- comb.samples(svy.A=svy.qA.h, svy.B=svy.qB.h,
  svy.C=svy.qC, y.lab="Lrn", z.lab="c.Days",
  form.x=~Eth:Sex:Age-1, estimation="synthetic",
  micro=TRUE, calfun="linear", bounds=c(.5,Inf),
  maxit=100)

head(out.4$Z.A)
head(out.4$Y.B)
```

 comp.cont

Compares two distributions of the same continuous variable

Description

This function estimate the “closeness” of the distributions of the same continuous variable(s) but estimated from different data sources.

Usage

```
comp.cont(data.A, data.B, xlab.A, xlab.B = NULL, w.A = NULL,
  w.B = NULL, ref = FALSE)
```

Arguments

| | |
|--------|--|
| data.A | A dataframe or matrix containing the variable of interest xlab.A and eventual associated survey weights w.A. |
| data.B | A dataframe or matrix containing the variable of interest xlab.B and eventual associated survey weights w.B. |
| xlab.A | Character string providing the name of the variable in data.A whose estimated distribution should be compared with that estimated from data.B. |
| xlab.B | Character string providing the name of the variable in data.B whose distribution should be compared with that estimated from data.A. If xlab.B=NULL (default) then it assumed xlab.B=xlab.A. |

| | |
|-----|--|
| w.A | Character string providing the name of the optional weighting variable in data.A that, in case, should be used to estimate the distribution of xlab.A |
| w.B | Character string providing the name of the optional weighting variable in data.B that, in case, should be used to estimate the distribution of xlab.B |
| ref | Logical. When ref = TRUE, the distribution of xlab.B estimated from data.B is considered the reference distribution (true or reliable estimate of distribution). |

Details

This function calculates well-known summary measures (min, Q1, median, mean, Q3, max and sd) estimated from the available data. It also compares the quantiles estimated from data.A with those estimated from data.B and returns the average of the absolute value of the differences and the average of the squared differences. Note that the number of percentiles estimated depends on the minimum between the two sample sizes. Note that the number of estimated percentiles depends on the minimum between the two sample sizes. Only quartiles are calculated if $\min(n.A, n.B) \leq 50$; quintiles are estimated if $\min(n.A, n.B) > 50$ and $\min(n.A, n.B) \leq 150$; deciles are estimated if $\min(n.A, n.B) > 150$ and $\min(n.A, n.B) \leq 250$; finally, quantiles for $\text{probs} = \text{seq}(\text{from} = 0.05, \text{to} = 0.95, \text{by} = 0.05)$ are estimated when $\min(n.A, n.B) > 250$. If survey weights are available (indicated by w.A and/or w.B), they are used to estimate the quantiles by calling the function [wtd.quantile](#) in the package **Hmisc**.

The dissimilarities between the estimated empirical distribution functions are calculated. The measures considered are the maximum value of the differences, the sum of the absolute values of the minimum and maximum, and the average of the absolute differences. If weights are given, they are used in the estimation of the empirical cumulative distribution function. Note that when ref=TRUE is given, the estimation of the density and the empirical cumulative distribution will be guided by the data in data.B.

Finally, the total variation distance, the overlap and the Hellinger are calculated on the transformed categorised variable. Note that the breaks to categorise the variable are decided according to the Freedman-Diaconis rule ([nclass](#)) and, in this case, with ref=TRUE the IQR is estimated on data.B alone, whereas with ref=FALSE it is estimated by combining the two data sources. If present, the weights are used to estimate the relative frequencies of the categorised variable. *total variation distance*:

$$\Delta_{AB} = \frac{1}{2} \sum_{j=1}^J |p_{A,j} - p_{B,j}|$$

where $p_{s,j}$ are the relative frequencies ($0 \leq p_{s,j} \leq 1$). The dissimilarity index ranges from 0 (minimum dissimilarity) to 1. The total variation distance comes along with its complement to 1, said “overlap” between distributions.

the *Hellinger’s distance*:

$$d_{H,AB} = \sqrt{\frac{1}{2} \sum_{j=1}^J (\sqrt{p_{A,j}} - \sqrt{p_{B,j}})^2}$$

It is a dissimilarity measure ranging from 0 (distributions are equal) to 1 (max dissimilarity). It satisfies all the properties of a distance measure ($0 \leq d_{H,AB} \leq 1$; symmetry and triangle inequality). Hellinger’s distance is related to the total variation distance, and it is possible to show that:

$$d_{H,AB}^2 \leq \Delta_{AB} \leq d_{H,AB}\sqrt{2}$$

Value

A list object with four components.

| | |
|------------|--|
| summary | A matrix with summaries of xlab.A estimated on data.A and summaries of xlab.B estimated on data.B |
| diff.Qs | Average of absolute and squared differences between the quantiles of xlab.A estimated on data.A and the corresponding ones of xlab.B estimated on data.B |
| dist.ecdf | Dissimilarity measures between the estimated empirical cumulative distribution functions. |
| dist.discr | Distance between the distributions after discretization of the target variable. |

Author(s)

Marcello D’Orazio <mdo.statmatch@gmail.com>

References

Bellhouse D.R. and J. E. Stafford (1999). “Density Estimation from Complex Surveys”. *Statistica Sinica*, **9**, 407–424.

See Also

[plotCont](#)

comp.prop

Compares two distributions of the same categorical variable

Description

This function compares two (estimated) distributions of the same categorical variable(s).

Usage

```
comp.prop(p1, p2, n1, n2=NULL, ref=FALSE)
```

Arguments

| | |
|----|--|
| p1 | A vector or an array containing relative or absolute frequencies for one or more categorical variables. Usually it is the output of the function <code>xtabs</code> or <code>table</code> . |
| p2 | A vector or an array containing relative or absolute frequencies for one or more categorical variables. Usually it is the output of the function <code>xtabs</code> or <code>table</code> . If <code>ref = FALSE</code> then p2 is a further estimate of the distribution of the categorical variable(s) being considered. On the contrary (<code>ref = TRUE</code>) it is the ‘reference’ distribution (the distribution considered true or a reliable estimate). |

| | |
|-----|--|
| n1 | The size of the sample on which p1 has been estimated. |
| n2 | The size of the sample on which p2 has been estimated, required just when ref = FALSE (p2 is estimated on another sample and is not the reference distribution). |
| ref | Logical. When ref = TRUE, p2 is the reference distribution (true or reliable estimate of distribution), on the contrary when ref = FALSE it an estimate of the distribution derived from another sample with sample size n2. |

Details

This function computes some similarity or dissimilarity measures between marginal (joint) distribution of categorical variables(s). The following measures are considered:

Dissimilarity index or total variation distance:

$$\Delta_{12} = \frac{1}{2} \sum_{j=1}^J |p_{1,j} - p_{2,j}|$$

where $p_{s,j}$ are the relative frequencies ($0 \leq p_{s,j} \leq 1$). The dissimilarity index ranges from 0 (minimum dissimilarity) to 1. It can be interpreted as the smallest fraction of units that need to be reclassified in order to make the distributions equal. When p2 is the reference distribution (true or expected distribution under a given hypothesis) than, following the Agresti's rule of thumb (Agresti 2002, pp. 329–330), values of $\Delta_{12} \leq 0.03$ denotes that the estimated distribution p1 follows the true or expected pattern quite closely.

Overlap between two distributions:

$$O_{12} = \sum_{j=1}^J \min(p_{1,j}, p_{2,j})$$

It is a measure of similarity which ranges from 0 to 1 (the distributions are equal). It is worth noting that $O_{12} = 1 - \Delta_{12}$.

Bhattacharyya coefficient:

$$B_{12} = \sum_{j=1}^J \sqrt{p_{1,j} \times p_{2,j}}$$

It is a measure of similarity and ranges from 0 to 1 (the distributions are equal).

Hellinger's distance:

$$d_{H,12} = \sqrt{1 - B_{12}}$$

It is a dissimilarity measure ranging from 0 (distributions are equal) to 1 (max dissimilarity). It satisfies all the properties of a distance measure ($0 \leq d_{H,12} \leq 1$; symmetry and triangle inequality). Hellinger's distance is related to the dissimilarity index, and it is possible to show that:

$$d_{H,12}^2 \leq \Delta_{12} \leq d_{H,12} \sqrt{2}$$

Alongside with those similarity/dissimilarity measures the Pearson's Chi-squared is computed. Two formulas are considered. When p_2 is the reference distribution (true or expected under some hypothesis, $ref=TRUE$):

$$\chi_P^2 = n_1 \sum_{j=1}^J \frac{(p_{1,j} - p_{2,j})^2}{p_{2,j}}$$

When p_2 is a distribution estimated on a second sample then:

$$\chi_P^2 = \sum_{i=1}^2 \sum_{j=1}^J n_i \frac{(p_{i,j} - p_{+,j})^2}{p_{+,j}}$$

where $p_{+,j}$ is the expected frequency for category j , obtained as follows:

$$p_{+,j} = \frac{n_1 p_{1,j} + n_2 p_{2,j}}{n_1 + n_2}$$

being n_1 and n_2 the sizes of the samples.

The Chi-Square value can be used to test the hypothesis that two distributions are equal ($df = J - 1$). Unfortunately such a test would not be useful when the distribution are estimated from samples selected from a finite population using complex selection schemes (stratification, clustering, etc.). In such a case different alternative corrected Chi-square tests are available (cf. Sarndal et al., 1992, Sec. 13.5). One possibility consist in dividing the Pearson's Chi-square test by the *generalised design effect* of both the surveys. Its estimation is not straightforward (sampling design variables need to be available). Generally speaking, the generalised design effect is smaller than 1 in the presence of stratified random sampling designs, while it exceeds 1 the presence of a two stage cluster sampling design. For the purposes of analysis it is reported the value of the generalised design effect g that would determine the acceptance of the null hypothesis (equality of distributions) in the case of $\alpha = 0.05$ ($df = J - 1$), i.e. values of g such that

$$\frac{\chi_P^2}{g} \leq \chi_{J-1,0.05}^2$$

Value

A list object with two or three components depending on the argument `ref`.

| | |
|---------------------|--|
| <code>meas</code> | A vector with the measures of similarity/dissimilarity between the distributions: dissimilarity index (" <code>tvd</code> "), overlap (" <code>overlap</code> "), Bhattacharyya coefficient (" <code>Bhatt</code> ") and Hellinger's distance (" <code>Hell</code> "). |
| <code>chi.sq</code> | A vector with the following values: Pearson's Chi-square (" <code>Pearson</code> "), the degrees of freedom (" <code>df</code> "), the percentile of a Chi-squared distribution (" <code>q0.05</code> ") and the largest admissible value of the generalised design effect that would determine the acceptance of H_0 (equality of distributions). |
| <code>p.exp</code> | When <code>ref=FALSE</code> it is reported the value of the reference distribution $p_{+,j}$ estimated used in deriving the Chi-square statistic and also the dissimilarity index. On the contrary (<code>ref=TRUE</code>) it is set equal to the argument <code>p2</code> . |

Author(s)

Marcello D’Orazio <mdo.statmatch@gmail.com>

References

Agresti A (2002) *Categorical Data Analysis. Second Edition*. Wiley, new York.
 Sarndal CE, Swensson B, Wretman JH (1992) *Model Assisted Survey Sampling*. Springer–Verlag, New York.

Examples

```
data(quine, package="MASS") #loads quine from MASS
str(quine)

# split quine in two subsets
suppressWarnings(RNGversion("3.5.0"))
set.seed(124)
lab.A <- sample(nrow(quine), 70, replace=TRUE)
quine.A <- quine[lab.A, c("Eth", "Sex", "Age")]
quine.B <- quine[-lab.A, c("Eth", "Sex", "Age")]

# compare est. distributions from 2 samples
# 1 variable
tt.A <- xtabs(~Age, data=quine.A)
tt.B <- xtabs(~Age, data=quine.B)
comp.prop(p1=tt.A, p2=tt.B, n1=nrow(quine.A), n2=nrow(quine.B), ref=FALSE)

# joint distr. of more variables
tt.A <- xtabs(~Eth+Sex+Age, data=quine.A)
tt.B <- xtabs(~Eth+Sex+Age, data=quine.B)
comp.prop(p1=tt.A, p2=tt.B, n1=nrow(quine.A), n2=nrow(quine.B), ref=FALSE)

# compare est. distr. with a one considered as reference
tt.A <- xtabs(~Eth+Sex+Age, data=quine.A)
tt.all <- xtabs(~Eth+Sex+Age, data=quine)
comp.prop(p1=tt.A, p2=tt.all, n1=nrow(quine.A), n2=NULL, ref=TRUE)
```

create.fused

Creates a matched (synthetic) dataset

Description

Creates a *synthetic* data frame after the statistical matching of two data sources at *micro* level.

Usage

```
create.fused(data.rec, data.don, mtc.ids,
             z.vars, dup.x=FALSE, match.vars=NULL)
```

Arguments

| | |
|------------|---|
| data.rec | A matrix or data frame that plays the role of <i>recipient</i> in the statistical matching application. |
| data.don | A matrix or data frame that plays the role of <i>donor</i> in the statistical matching application. |
| mtc.ids | A matrix with two columns. Each row must contain the name or the index of the recipient record (row) in data.don and the name or the index of the corresponding donor record (row) in data.don. Note that this type of matrix is returned by the functions NND.hotdeck , RANDwNND.hotdeck , rankNND.hotdeck , and mixed.mtc . |
| z.vars | A character vector with the names of the variables available only in data.don that should be “donated” to data.rec. |
| dup.x | Logical. When TRUE the values of the matching variables in data.don are also “donated” to data.rec. The names of the matching variables have to be specified with the argument match.vars. To avoid confusion, the matching variables added to data.rec are renamed by adding the suffix “don”. By default dup.x=FALSE. |
| match.vars | A character vector with the names of the matching variables. It has to be specified only when dup.x=TRUE. |

Details

This function allows to create the synthetic (or fused) data set after the application of a statistical matching in a *micro* framework. For details see D’Orazio *et al.* (2006).

Value

The data frame data.rec with the z.vars filled in and, when dup.x=TRUE, with the values of the matching variables match.vars observed on the donor records.

Author(s)

Marcello D’Orazio <mdo.statmatch@gmail.com>

References

D’Orazio, M., Di Zio, M. and Scanu, M. (2006). *Statistical Matching: Theory and Practice*. Wiley, Chichester.

See Also

[NND.hotdeck](#) [RANDwNND.hotdeck](#) [rankNND.hotdeck](#)

Examples

```

lab <- c(1:15, 51:65, 101:115)
iris.rec <- iris[lab, c(1:3,5)] # recipient data.frame
iris.don <- iris[-lab, c(1:2,4:5)] # donor data.frame

# Now iris.rec and iris.don have the variables
# "Sepal.Length", "Sepal.Width" and "Species"
# in common.
# "Petal.Length" is available only in iris.rec
# "Petal.Width" is available only in iris.don

# find the closest donors using NND hot deck;
# distances are computed on "Sepal.Length" and "Sepal.Width"

out.NND <- NND.hotdeck(data.rec=iris.rec, data.don=iris.don,
  match.vars=c("Sepal.Length", "Sepal.Width"),
  don.class="Species")

# create synthetic data.set, without the
# duplication of the matching variables

fused.0 <- create.fused(data.rec=iris.rec, data.don=iris.don,
  mtc.ids=out.NND$mtc.ids, z.vars="Petal.Width")

# create synthetic data.set, with the "duplication"
# of the matching variables

fused.1 <- create.fused(data.rec=iris.rec, data.don=iris.don,
  mtc.ids=out.NND$mtc.ids, z.vars="Petal.Width",
  dup.x=TRUE, match.vars=c("Sepal.Length", "Sepal.Width"))

```

| | |
|----------------|--|
| create.imputed | <i>Fills-in missing values in the recipient dataset with values observed on the donors units</i> |
|----------------|--|

Description

Imputes the missing values (NAs) in the recipient dataset with values observed on the donors units after search of donors with NND or random hotdeck.

Usage

```
create.imputed(data.rec, data.don, mtc.ids)
```

Arguments

| | |
|----------|--|
| data.rec | A matrix or data frame that has missing values. |
| data.don | A matrix or data frame that is used for donation (imputation). |

`mtc.ids` A matrix with two columns. Each row must contain the name or the index of the recipient record (row) in `data.don` and the name or the index of the corresponding donor record (row) in `data.don`. Note that this type of matrix is returned by the functions `NND.hotdeck`, `RANDwNND.hotdeck`, `rankNND.hotdeck`, and `mixed.mtc`.

Details

This function allows the missing values (NAs) in the recipient to be filled in with values observed in the donor dataset after searching for donors via NND or random hotdeck using the functions available in the package, i.e. `NND.hotdeck`, `RANDwNND.hotdeck`, `rankNND.hotdeck`, and `mixed.mtc`. If the same record in the recipient dataset has 2 or more NAs, they will all be replaced with the values observed for that unit on the selected donor; this is equivalent to joint hotdeck imputation.

Value

The data frame `data.rec` missing values (NAs) filled in.

Author(s)

Marcello D'Orazio <mdo.statmatch@gmail.com>

References

D'Orazio, M., Di Zio, M. and Scanu, M. (2006). *Statistical Matching: Theory and Practice*. Wiley, Chichester.

See Also

[NND.hotdeck](#) [RANDwNND.hotdeck](#) [rankNND.hotdeck](#)

Examples

```
# introduce missing values
# in Petal.Length of iris dataset

set.seed(13579)
pos <- sample(x = 1:nrow(iris), size = 15,
             replace = FALSE)

iris.rec <- iris[pos, ] # recipient data.frame with missing values
iris.rec[, "Petal.Length"] <- NA

iris.don <- iris[-pos, ] # donor data.frame ALL observed

# find the closest donors using NND hot deck;
# distances are computed on "Petal.Width"
# donors only of the same Specie

out.NND <- NND.hotdeck(data.rec=iris.rec, data.don=iris.don,
```



```
match.vars=c("Petal.Width"),
don.class="Species")

# impute missing
iris.rec.imp <- create.imputed(data.rec=iris.rec, data.don=iris.don,
                             mtc.ids=out.NND$mtc.ids)
summary(iris.rec.imp$Petal.Length)
```

fact2dummy

Transforms a categorical variable in a set of dummy variables

Description

Transforms a factor or more factors contained in a data frame in a set of dummy variables, while numeric variables remain unchanged.

Usage

```
fact2dummy(data, all=TRUE, lab="x")
```

Arguments

| | |
|------|---|
| data | A factor or a data frame that contains one or more factors (columns whose class is “factor” or “ordered”) that have to be substituted by the corresponding dummy variables. |
| all | Logical. When all=TRUE (default) the output matrix will contain as many dummy variables as the number of the levels of the factor variable. On the contrary, when all=FALSE, the dummy variable related to the last level of the factor is dropped. |
| lab | A character string with the name of the variable to be pasted with its levels. This is used only when data is a factor. By default it is set to “x”. |

Details

This function substitutes categorical variables in the input data frame (columns whose class is “factor” or “ordered”) with the corresponding dummy variables. Note that if a factor includes a missing values (NA) then all the associated dummies will report an NA in correspondence of the missing observation (row).

Value

A matrix with the dummy variables instead of initial factor variables.

Author(s)

Marcello D’Orazio <mdo.statmatch@gmail.com>

See Also[gower.dist](#)**Examples**

```
x <- runif(5)
y <- factor(c(1,2,1,2,2))
z <- ordered(c(1,2,3,2,2))
xyz <- data.frame(x,y,z)
fact2dummy(xyz)

fact2dummy(xyz, all=FALSE)

#example with iris data frame
str(iris)
ir.mat <- fact2dummy(iris)
head(ir.mat)
```

Fbounds.pred

Estimates Frechet bounds for cells in the contingency table crossing two categorical variables observed in distinct samples referred to the same target population.

Description

This function assesses the uncertainty in estimating the contingency table crossing `y.rec` (Y) and `z.don` (Z) when the two variables are observed in two different samples sharing a number of common predictors.

Usage

```
Fbounds.pred(data.rec, data.don,
             match.vars, y.rec, z.don, pred = "multinom",
             w.rec = NULL, w.don = NULL, type.pred = "random",
             out.pred = FALSE, ...)
```

Arguments

| | |
|-------------------------|--|
| <code>data.rec</code> | dataframe including the Xs (predictors, listed in <code>match.vars</code>) and <code>y.rec</code> (response; target variable in this dataset) |
| <code>data.don</code> | dataframe including the Xs (predictors, listed in <code>match.vars</code>) and <code>z.don</code> (response; target variable in this dataset) |
| <code>match.vars</code> | vector with the names of the Xs variables to be used as predictors (or set in which select the best predictors with lasso) of respectively <code>y.rec</code> and <code>z.don</code> |

| | |
|-----------|--|
| y.rec | character indicating the name of Y target variable in data.rec. It should be a factor. |
| z.don | character indicating the name of Z target variable in data.don. It should be a factor. |
| pred | character specifying the method used to obtain predictions of both Y and Z. Available methods include pred = "multinom" (default) fits two multinomial models (nnet function <code>multinom</code>) to get predictions with Y and Z as response variables and match.vars as predictors; pred = "lasso" uses the lasso method (R package glmnet , function <code>cv.glmnet</code>) and cross-validation to select a subset of match.vars that are the best predictors of Y and Z, respectively, and then fits the multinomial models with the selected predictors; pred = "nb" uses the Naive Bayes classifier to get predictions of Y and Z respectively (R package naivebayes function <code>naive_bayes</code>); pred = "rf" fits randomForest to get predictions of both Y and Z (function <code>randomForest</code> in randomForest). |
| w.rec | name of the variable with the weights of the units in data.rec, if available (default is NULL); the weights, if available, are only used for estimating bounds, not for fitting models. |
| w.don | name of the variable with the weights of the units in data.don, if available (default is NULL); the weights, if available, are only used for estimating bounds, not for fitting models. |
| type.pred | string specifying how to obtain the predictions of Y and Z. By default, the fitted models return conditional probabilities (or scores), then if type.pred = "random" (default), the predicted class of Y (Z) is obtained by a random draw with selection probabilities equal to the estimated conditional probabilities (scores); on the contrary, if type.pred = "mostvoted", the predicted class is the one with the highest estimated conditional probability (score). |
| out.pred | Logical. If TRUE (default is FALSE) returns the input datasets with the estimated conditional probabilities (depending on pred argument), the predicted class for the target variable (Y or Z) in the dataset (depending on type.pred argument) and the true observed class of Y (or Z). |
| ... | additional arguments, if needed. |

Details

The function evaluates the uncertainty in estimating the contingency table crossing y.rec (Y) and z.don (Z) when the two variables are observed in two different samples related to the same target population, but the samples share a number of common predictors. The evaluation of the uncertainty is equivalent to estimating the bounds for each cell in the contingency table where Y and Z intersect; the bounds can be unconditional (Frechet property) or conditional on the predictions of both Y and Z provided by the models fitted according to the pred argument. This latter way of working avoids many of the drawbacks of obtaining expectations of conditional bounds when conditioning on many X variables, and allows the inclusion of non-categorical predictors. The final estimation of the conditional bounds is provided by the function `Frechet.bounds.cat`.

Value

a list with the following components:

up.rec only when out.pred = TRUE it corresponds to a smaller version of data.rec with the estimated conditional probabilities for both Y and Z (depending on pred argument), the predicted class of Y (depending on type.pred argument), the predicted class of Z (depending on type.pred argument), the true observed class of Y and the predictors (argument match.vars) (and the weights when w.rec is specified).

up.don only when out.pred = TRUE it corresponds to a smaller version of data.don with the estimated conditional probabilities for both Y and Z (depending on pred argument), the predicted class of Y (depending on type.pred argument), the predicted class of Z (depending on type.pred argument), the true observed class of Z and the predictors (argument match.vars) (and the weights when w.don is specified).

p.xx.ini the estimated relative frequencies in the table crossing predictions of Y and Z; it is estimated after pooling the samples (weighted average of estimates obtained on the separates samples);

p.xy.ini the estimated table crossing Y and the predictions of both Y and Z estimated from data.rec (weights are used if provided with the w.rec argument);

p.xz.ini the estimated table crossing Z and the predictions of both Y and Z estimated from data.don (weights are used if provided with the w.don argument);

accuracy the estimated accuracy in predicting respectively Y and Z with the chosen method (argument pred) and the available predictors (argument match.vars);

bounds a data.frame whose columns reports the estimated unconditional and conditional bounds for each cell in the contingency table crossing y.rec(Y) and z.don(Z);

uncertainty the uncertainty associated to input data, measured in terms of average width of uncertainty bounds with and without conditioning on the predictions (for further details see [Frechet.bounds.cat](#)).

Author(s)

Marcello D’Orazio <mdo.statmatch@gmail.com>

References

D’Orazio, M., (2024). *Is Statistical Matching feasible?* Note, https://www.researchgate.net/publication/387699016_Is_statistical_matching_feasible.

See Also

[Frechet.bounds.cat](#).

Examples

```
data(quine, package="MASS") #loads quine from MASS
str(quine)

# split quine in two subsets
set.seed(223344)
lab.A <- sample(nrow(quine), 70, replace=TRUE)
quine.A <- quine[lab.A, 1:3]
quine.B <- quine[-lab.A, 2:4]

# multinomial model and predictions with most-voted criterion
```

```

fbp <- Fbounds.pred(data.rec = quine.A, data.don = quine.B,
  match.vars = c("Sex", "Age"),
  y.rec = "Eth", z.don = "Lrn",
  pred = "multinom", type.pred = "mostvoted")

fbp$p.xx.ini # estimated cross-tab of predictions
fbp$bounds # estimated conditional and unconditional bounds
fbp$uncertainty # estimated uncertainty about Y*Z

# multinomial model and predictions with randomized criterion
fbp <- Fbounds.pred(data.rec = quine.A, data.don = quine.B,
  match.vars = c("Sex", "Age"),
  y.rec = "Eth", z.don = "Lrn",
  pred = "multinom", type.pred = "random")

fbp$p.xx.ini # estimated cross-tab of predictions
fbp$bounds # estimated conditional and unconditional bounds
fbp$uncertainty # estimated uncertainty about Y*Z

```

| | |
|---------------|---|
| Fbwidths.by.x | <i>Computes the Frechet bounds of cells in a contingency table by considering all the possible subsets of the common variables.</i> |
|---------------|---|

Description

This function permits to compute the bounds for cell probabilities in the contingency table Y vs. Z starting from the marginal tables (X vs. Y), (X vs. Z) and the joint distribution of the X variables, by considering all the possible subsets of the X variables. In this manner it is possible to identify which subset of the X variables produces the major reduction of the average width of conditional bounds.

Usage

```
Fbwidths.by.x(tab.x, tab.xy, tab.xz, deal.sparse="discard",
  nA=NULL, nB=NULL, ...)
```

Arguments

| | |
|--------|---|
| tab.x | A R table crossing the X variables. This table must be obtained by using the function <code>xtabs</code> or <code>table</code> , e.g. <code>tab.x <- xtabs(~x1+x2+x3, data=data.all)</code> . |
| tab.xy | A R table of X vs. Y variable. This table must be obtained by using the function <code>xtabs</code> or <code>table</code> , e.g. <code>table.xy <- xtabs(~x1+x2+x3+y, data=data.A)</code> . A single categorical Y variables is allowed. One or more categorical variables can be considered as X variables (common variables). The same X variables in <code>tab.x</code> must be available in <code>tab.xy</code> . Moreover, it is assumed that the joint distribution of the X variables computed from <code>tab.xy</code> is equal to <code>tab.x</code> ; a warning is produced if this is not true. |

| | |
|-------------|--|
| tab.xz | <p>A R table of X vs. Z variable. This table must be obtained by using the function <code>xtabs</code> or <code>table</code>, e.g.</p> <pre>tab.xz <- xtabs(~x1+x2+x3+z, data=data.B).</pre> <p>A single categorical Z variable is allowed. One or more categorical variables can be considered as X variables (common variables). The same X variables in <code>tab.x</code> must be available in <code>tab.xz</code>. Moreover, it is assumed that the joint distribution of the X variables computed from <code>tab.xz</code> is equal to <code>tab.x</code>; a warning is produced if this is not true.</p> |
| deal.sparse | <p>Text, how to estimate the cell relative frequencies when dealing with too sparse tables. When <code>deal.sparse="discard"</code> (default) no estimation is performed if <code>tab.xy</code> or <code>tab.xz</code> is too sparse. When <code>deal.sparse="relfreq"</code> the standard estimator (cell count divided by the sample size) is considered. Note that here sparseness is measured by number of cells with respect to the sample size; sparse table are those where the number of cells exceeds the sample size (see Details).</p> |
| nA | <p>Integer, sample size of file A used to estimate <code>tab.xy</code>. If NULL, it is obtained as sum of frequencies in <code>tab.xy</code>.</p> |
| nB | <p>Integer, sample size of file B used to estimate <code>tab.xz</code>. If NULL, it is obtained as sum of frequencies in <code>tab.xz</code>.</p> |
| ... | <p>Additional arguments that may be required when deriving an estimate of uncertainty by calling <code>Frechet.bounds.cat</code>.</p> |

Details

This function permits to compute the Frechet bounds for the frequencies in the contingency table of **Y** vs. **Z**, starting from the conditional distributions $P(Y|X)$ and $P(Z|X)$ (for details see `Frechet.bounds.cat`), by considering all the possible subsets of the **X** variables. In this manner it is possible to identify the subset of the **X** variables, with highest association with both **Y** and **Z**, that permits to reduce the uncertainty concerning the distribution of **Y** vs. **Z**.

The uncertainty is measured by the average of the widths of the bounds for the cells in the table **Y** vs. **Z**:

$$\bar{d} = \frac{1}{J \times K} \sum_{j,k} (p_{Y=j,Z=k}^{(up)} - p_{Y=j,Z=k}^{(low)})$$

For details see `Frechet.bounds.cat`.

Provided that uncertainty, measured in terms of \bar{d} , tends to reduce when conditioning on a higher number of **X** variables. Two penalties are introduced to account for the additional number of cells to be estimated when adding a **X** variable. The first penalty, introduced in D'Orazio et al. (2017), is:

$$g_1 = \log \left(1 + \frac{H_{D_m}}{H_{D_Q}} \right)$$

Where H_{D_m} is the number of cell in the table obtained by crossing the given subset of **X** variables and the H_{D_Q} is the number of cell in the table achieved by crossing all the available **X** variables. A second penalty takes into account the number of cells to estimate with respect to the sample size (D'Orazio et al., 2019). It is obtained as:

$$g_2 = \max \left[\frac{1}{n_A - H_{D_m} \times J}, \frac{1}{n_B - H_{D_m} \times K} \right]$$

with $n_A > H_{D_m} \times J$ and $n_B > H_{D_m} \times K$. In practice, it is considered the number of cells to estimate compared to the sample size. This criterion is considered to measure sparseness too. In particular, for the purposes of this function, tables are NOT considered sparse when:

$$\min \left[\frac{n_A}{H_{D_m} \times J}, \frac{n_B}{H_{D_m} \times K} \right] > 1$$

This rule is applied when deciding how to proceed with estimation in case of sparse table (argument `deal.sparse`). Note that sparseness can be measured in different manners. The outputs include also the empty cells in each table (due to statistical zeros or structural zeros) and the Cohen's effect size with respect to the case of uniform distribution of frequencies across cells (the value $1/\text{no.of.cells}$ in every cell):

$$\omega_{eq} = \sqrt{H \sum_{h=1}^H (\hat{p}_h - 1/H)^2}$$

values of ω_{eq} jointly with $n/H \leq 1$ usually indicate severe sparseness.

Value

A list with the estimated bounds for the cells in the table of Y vs. Z for each possible subset of the X variables. The final component in the list, `sum.unc`, is a data.frame that summarizes the main results. In particular, it reports the number of X variables ("`x.vars`"), the number of cells in each of the input tables and the cells with frequency equal to 0 (columns ending with `freq0`). Moreover, it reported the value ("`av.n`") of the rule used to decide whether we are dealing with a sparse case (see Details) and the Cohen's effect size measured for the table crossing the considered combination of the X variables. Finally, it is provided the average width of the uncertainty intervals ("`av.width`"), the penalty terms `g1` and `g2` ("`penalty1`" and "`penalty2`" respectively), and the penalized average widths ("`av.width.pen1`" and "`av.width.pen2`", where `av.width.pen1=av.width+pen1` and `av.width.pen2=av.width+pen2`).

Author(s)

Marcello D'Orazio <`mdo.statmatch@gmail.com`>

References

- D'Orazio, M., Di Zio, M. and Scanu, M. (2006). *Statistical Matching: Theory and Practice*. Wiley, Chichester.
- D'Orazio, M., Di Zio, M. and Scanu, M. (2017). "The use of uncertainty to choose matching variables in statistical matching". *International Journal of Approximate Reasoning*, 90, pp. 433-440.
- D'Orazio, M., Di Zio, M. and Scanu, M. (2019). "Auxiliary variable selection in a a statistical matching problem". In Zhang, L.-C. and Chambers, R. L. (eds.) *Analysis of Integrated Data*, Chapman & Hall/CRC (Forthcoming).

See Also

[Frechet.bounds.cat](#), [harmonize.x](#)

Examples

```
# un-comment to run
#
# data(quine, package="MASS") #loads quine from MASS
# str(quine)
# quine$c.Days <- cut(quine$Days, c(-1, seq(0,50,10),100))
# table(quine$c.Days)
#
#
# # split quine in two subsets
# suppressWarnings(RNGversion("3.5.0"))
# set.seed(4567)
# lab.A <- sample(nrow(quine), 70, replace=TRUE)
# quine.A <- quine[lab.A, 1:4]
# quine.B <- quine[-lab.A, c(1:3,6)]
#
# # compute the tables required by Fbwidths.by.x()
# freq.xA <- xtabs(~Eth+Sex+Age, data=quine.A)
# freq.xB <- xtabs(~Eth+Sex+Age, data=quine.B)
#
# freq.xy <- xtabs(~Eth+Sex+Age+Lrn, data=quine.A)
# freq.xz <- xtabs(~Eth+Sex+Age+c.Days, data=quine.B)
#
# # apply Fbwidths.by.x()
# bounds.yz <- Fbwidths.by.x(tab.x=freq.xA+freq.xB, tab.xy=freq.xy,
#                             tab.xz=freq.xz)
#
# bounds.yz$sum.unc
```

Frechet.bounds.cat *Frechet bounds of cells in a contingency table*

Description

This function permits to derive the bounds for cell probabilities of the table Y vs. Z starting from the marginal tables (X vs. Y), (X vs. Z) and the joint distribution of the X variables.

Usage

```
Frechet.bounds.cat(tab.x, tab.xy, tab.xz, print.f="tables", align.margins = FALSE,
                  tol= 0.001, warn = TRUE)
```


Arguments

| | |
|---------------|--|
| tab.x | <p>A R table crossing the X variables. This table must be obtained by using the function <code>xtabs</code> or <code>table</code>, e.g.</p> <pre>tab.x <- xtabs(~x1+x2+x3, data=data.all).</pre> <p>When <code>tab.x = NULL</code> then only <code>tab.xy</code> and <code>tab.xz</code> must be supplied.</p> |
| tab.xy | <p>A R table of X vs. Y variable. This table must be obtained by using the function <code>xtabs</code> or <code>table</code>, e.g.</p> <pre>table.xy <- xtabs(~x1+x2+x3+y, data=data.A).</pre> <p>A single categorical Y variable is allowed. One or more categorical variables can be considered as X variables (common variables). Obviously, the same X variables in <code>tab.x</code> must be available in <code>tab.xy</code>. Usually, it is assumed that the joint distribution of the X variables computed from <code>tab.xy</code> is equal to <code>tab.x</code> (a warning appears if any absolute difference is greater than <code>tol</code>). Note that when marginal distribution of X in <code>tab.xy</code> is not equal to that of <code>tab.x</code> it is possible to ask their alignment (see argument <code>align.margins</code>).</p> <p>When <code>tab.x = NULL</code> then <code>tab.xy</code> should be a one-dimensional table providing the marginal distribution of the Y variable.</p> |
| tab.xz | <p>A R table of X vs. Z variable. This table must be obtained by using the function <code>xtabs</code> or <code>table</code>, e.g.</p> <pre>tab.xz <- xtabs(~x1+x2+x3+z, data=data.B).</pre> <p>A single categorical Z variable is allowed. One or more categorical variables can be considered as X variables (common variables). The same X variables in <code>tab.x</code> must be available in <code>tab.xz</code>. Usually, it is assumed that the joint distribution of the X variables computed from <code>tab.xz</code> is equal to <code>tab.x</code> (a warning appears if any absolute difference is greater than <code>tol</code>). Note that when marginal distribution of X in <code>tab.xz</code> is not equal to that of <code>tab.x</code> it is possible to ask their alignment (see argument <code>align.margins</code>).</p> <p>When <code>tab.x = NULL</code> then <code>tab.xz</code> should be a one-dimensional table providing the marginal distribution of the Z variable.</p> |
| print.f | <p>A string: when <code>print.f="tables"</code> (default) all the cells' estimates will be saved as tables in a list. On the contrary, if <code>print.f="data.frame"</code>, they will be saved as columns of a data.frame.</p> |
| align.margins | <p>Logical (default FALSE). When when TRUE the distribution of X variables in <code>tab.xy</code> is aligned with the distribution resulting from <code>tab.x</code>, without affecting the marginal distribution of Y. Similarly, the distribution of X variables in <code>tab.xz</code> is aligned with the distribution resulting from <code>tab.x</code> without affecting the marginal distribution of Z. The alignment is performed by running IPF algorithm as implemented in the function <code>Estimate</code> in the package <code>mipfp</code>. Note that to avoid lack of convergence due to combinations of Xs encountered in one table but not in the other (statistical 0s), before running IPF a small constant ($1e-06$) is added to empty cells in <code>tab.xy</code> and <code>tab.xz</code>.</p> |
| tol | <p>Tolerance used in comparing joint distributions as far as X variables are considered (default <code>tol=0.001</code>); estimation of cells bounds would require that distribution of X variables computed from <code>tab.xy</code> and <code>tab.xz</code> should be approximately equal to that in <code>tab.x</code>, on contrary incoherences in estimated cells' bounds could happen. In case of not-coherent marginal distributions it is suggested to get them aligned by setting <code>align.margins=TRUE</code>.</p> |

warn Logical, when TRUE (default) return warnings when marginal distributions of **X** variables show differences greater than tol.

Details

This function permits to compute the expected conditional Frechet bounds for the relative frequencies in the contingency table of **Y** vs. **Z**, starting from the distributions $P(Y|X)$, $P(Z|X)$ and $P(X)$. The expected conditional bounds for the relative frequencies $p_{j,k}$ in the table **Y** vs. **Z** are:

$$p_{Y=j,Z=k}^{(low)} = \sum_i p_{X=i} \max(0; p_{Y=j|X=i} + p_{Z=k|X=i} - 1)$$

$$p_{Y=j,Z=k}^{(up)} = \sum_i p_{X=i} \min(p_{Y=j|X=i}; p_{Z=k|X=i})$$

The relative frequencies $p_{X=i} = n_i/n$ are computed from the frequencies in `tab.x`; the relative frequencies $p_{Y=j|X=i} = n_{ij}/n_{i+}$ are derived from `tab.xy`, finally, $p_{Z=k|X=i} = n_{ik}/n_{i+}$ are derived from `tab.xz`.

Estimation requires that all the starting tables share the same marginal distribution of the **X** variables.

This function returns also the unconditional bounds for the relative frequencies in the contingency table of **Y** vs. **Z**, i.e. computed also without considering the **X** variables:

$$\max\{0; p_{Y=j} + p_{Z=k} - 1\} \leq p_{Y=j,Z=k} \leq \min\{p_{Y=j}; p_{Z=k}\}$$

These bounds represent the unique output when `tab.x = NULL`.

Finally, the contingency table of **Y** vs. **Z** estimated under the Conditional Independence Assumption (CIA) is obtained by considering:

$$p_{Y=j,Z=k} = p_{Y=j|X=i} \times p_{Z=k|X=i} \times p_{X=i}.$$

When `tab.x = NULL` then it is also provided the expected table under the assumption of independence between **Y** and **Z**:

$$p_{Y=j,Z=k} = p_{Y=j} \times p_{Z=k}.$$

The presence of too many cells with 0s in the input contingency tables is an indication of sparseness; this is an unappealing situation when estimating the cells' relative frequencies needed to derive the bounds; in such cases the corresponding results may be unreliable. A possible alternative way of working consists in estimating the required parameters by considering a pseudo-Bayes estimator (see [pBayes](#)); in practice the input `tab.x`, `tab.xy` and `tab.xz` should be the ones provided by the [pBayes](#) function.

Value

When `print.f="tables"` (default) a list with the following components:

| | |
|--------------------------|---|
| <code>low.u</code> | The estimated lower bounds for the relative frequencies in the table Y vs. Z without conditioning on the X variables. |
| <code>up.u</code> | The estimated upper bounds for the relative frequencies in the table Y vs. Z without conditioning on the X variables. |
| <code>CIA</code> | The estimated relative frequencies in the table Y vs. Z under the Conditional Independence Assumption (CIA). |
| <code>low.cx</code> | The estimated lower bounds for the relative frequencies in the table Y vs. Z when conditioning on the X variables. |
| <code>up.cx</code> | The estimated upper bounds for the relative frequencies in the table Y vs. Z when conditioning on the X variables. |
| <code>uncertainty</code> | The uncertainty associated to input data, measured in terms of average width of uncertainty bounds with and without conditioning on the X variables. |

When `print.f="data.frame"` the output list contains just two components:

| | |
|--------------------------|---|
| <code>bounds</code> | A data.frame whose columns reports the estimated uncertainty bounds. |
| <code>uncertainty</code> | The uncertainty associated to input data, measured in terms of average width of uncertainty bounds with and without conditioning on the X variables. |

Author(s)

Marcello D’Orazio <mdo.statmatch@gmail.com>

References

- D’Orazio, M., Di Zio, M. and Scanu, M. (2006) “Statistical Matching for Categorical Data: Displaying Uncertainty and Using Logical Constraints”, *Journal of Official Statistics*, 22, pp. 137–157.
- D’Orazio, M., Di Zio, M. and Scanu, M. (2006). *Statistical Matching: Theory and Practice*. Wiley, Chichester.

See Also

[Fbwidths.by.x](#), [harmonize.x](#)

Examples

```
data(quine, package="MASS") #loads quine from MASS
str(quine)

# split quine in two subsets
suppressWarnings(RNGversion("3.5.0"))
set.seed(7654)
lab.A <- sample(nrow(quine), 70, replace=TRUE)
quine.A <- quine[lab.A, 1:3]
quine.B <- quine[-lab.A, 2:4]
```

```

# compute the tables required by Frechet.bounds.cat()
freq.xA <- xtabs(~Sex+Age, data=quine.A)
freq.xB <- xtabs(~Sex+Age, data=quine.B)

freq.xy <- xtabs(~Sex+Age+Eth, data=quine.A)
freq.xz <- xtabs(~Sex+Age+Lrn, data=quine.B)

# apply Frechet.bounds.cat()
bounds.yz <- Frechet.bounds.cat(tab.x=freq.xA+freq.xB, tab.xy=freq.xy,
                               tab.xz=freq.xz, print.f="data.frame")
bounds.yz

# harmonize distr. of Sex vs. Age during computations
# in Frechet.bounds.cat()

#compare marg. distribution of Xs in A and B vs. pooled estimate
comp.prop(p1=margin.table(freq.xy,c(1,2)), p2=freq.xA+freq.xB,
          n1=nrow(quine.A), n2=nrow(quine.A)+nrow(quine.B), ref=TRUE)

comp.prop(p1=margin.table(freq.xz,c(1,2)), p2=freq.xA+freq.xB,
          n1=nrow(quine.A), n2=nrow(quine.A)+nrow(quine.B), ref=TRUE)

bounds.yz <- Frechet.bounds.cat(tab.x=freq.xA+freq.xB, tab.xy=freq.xy,
                               tab.xz=freq.xz, print.f="data.frame", align.margins=TRUE)
bounds.yz

```

gower.dist

Computes the Gower's Distance

Description

This function computes the Gower's distance (dissimilarity) between units in a dataset or between observations in two distinct datasets.

Usage

```
gower.dist(data.x, data.y=data.x, rngs=NULL, KR.corr=TRUE, var.weights = NULL, robcb=NULL)
```

Arguments

`data.x` A matrix or a data frame containing variables that should be used in the computation of the distance.

Columns of mode numeric will be considered as interval scaled variables; columns of mode character or class factor will be considered as categorical nominal variables; columns of class ordered will be considered as categorical ordinal variables and, columns of mode logical will be considered as binary asymmetric variables (see Details for further information).

| | |
|-------------|---|
| | Missing values (NA) are allowed. If only data.x is supplied, the dissimilarities between rows of data.x will be computed. |
| data.y | A numeric matrix or data frame with the same variables, of the same type, as those in data.x. Dissimilarities between rows of data.x and rows of data.y will be computed. If not provided, by default it is assumed equal to data.x and only dissimilarities between rows of data.x will be computed. |
| rngs | A vector with the ranges to scale the variables. Its length must be equal to number of variables in data.x. In correspondence of non-numeric variables, just put 1 or NA. When rngs=NULL (default) the range of a numeric variable is estimated by jointly considering the values for the variable in data.x and those in data.y. Therefore, assuming rngs=NULL, if a variable "X1" is considered: <pre>rngs["X1"] <- max(data.x[, "X1"], data.y[, "X1"]) - min(data.x[, "X1"], data.y[, "X1"])</pre> |
| KR.corr | When TRUE (default) the extension of the Gower's dissimilarity measure proposed by Kaufman and Rousseeuw (1990) is used. Otherwise, when KR.corr=FALSE, the Gower's (1971) formula is considered. |
| var.weights | By default (NULL) each variable has the same weight (value 1) when calculating the overall distance (weighted average of distances on single variables; see Details). User can specify different weights for the different variables by providing a numeric value for each of the variables contributing to the distance. In other words, var.weights should be set equal to a numeric vector having length equal to the number of variables considered in calculating distance. Entered weights are scales to sum up to 1. |
| robcb | By default is (NULL). If robcb="boxp" the scaling of the Manhattan distance is done by using the difference between upper and lower fences of the Boxplot with k=3. In alternative, robcb="asyboxp" the scaling of the Manhattan distance is done by the difference between upper and lower fences of the modified Boxplot to account for slight skewness. In this case scaled distances greater than 1 are set equal to 1. This option is suggested in the presence of outliers in the continuous variables. |

Details

This function computes distances between records when variables of different type (categorical and continuous) have been observed. In order to handle different types of variables, the Gower's dissimilarity coefficient (Gower, 1971) is used. By default (KR.corr=TRUE) the Kaufman and Rousseeuw (1990) extension of the Gower's dissimilarity coefficient is used.

The final dissimilarity between the i th and j th unit is obtained as a weighted sum of dissimilarities for each variable:

$$d(i, j) = \frac{\sum_k \delta_{ijk} d_{ijk} w_k}{\sum_k \delta_{ijk} w_k}$$

In particular, d_{ijk} represents the distance between the i th and j th unit computed considering the k th variable, while w_k is the weight assigned to variable k (by default 1 for all the variables, unless

different weights are provided by user with argument `var.weights`). Distance depends on the nature of the variable:

- logical columns are considered as asymmetric binary variables, for such case $d_{ijk} = 0$ if $x_{ik} = x_{jk} = \text{TRUE}$, 1 otherwise;
- factor or character columns are considered as categorical nominal variables and $d_{ijk} = 0$ if $x_{ik} = x_{jk}$, 1 otherwise;
- numeric columns are considered as interval-scaled variables and

$$d_{ijk} = \frac{|x_{ik} - x_{jk}|}{R_k}$$

being R_k the range of the k th variable. The range is the one supplied with the argument `rngs` (`rngs[k]`) or the one computed on available data (when `rngs=NULL`);

- ordered columns are considered as categorical ordinal variables and the values are substituted with the corresponding position index, r_{ik} in the factor levels. When `KR.corr=FALSE` these position indexes (that are different from the output of the R function `rank`) are transformed in the following manner

$$z_{ik} = \frac{(r_{ik} - 1)}{\max(r_{ik}) - 1}$$

These new values, z_{ik} , are treated as observations of an interval scaled variable.

As far as the weight δ_{ijk} is concerned:

- $\delta_{ijk} = 0$ if $x_{ik} = \text{NA}$ or $x_{jk} = \text{NA}$;
- $\delta_{ijk} = 0$ if the variable is asymmetric binary and $x_{ik} = x_{jk} = 0$ or $x_{ik} = x_{jk} = \text{FALSE}$;
- $\delta_{ijk} = 1$ in all the other cases.

In practice, NAs and couple of cases with $x_{ik} = x_{jk} = \text{FALSE}$ do not contribute to distance computation.

Value

A matrix object with distances between rows of data `x` and those of data `y`.

Author(s)

Marcello D’Orazio <mdo.statmatch@gmail.com>

References

Gower, J. C. (1971), “A general coefficient of similarity and some of its properties”. *Biometrics*, **27**, 623–637.

Kaufman, L. and Rousseeuw, P.J. (1990), *Finding Groups in Data: An Introduction to Cluster Analysis*. Wiley, New York.

See Also

[daisy](#), [dist](#)

Examples

```
x1 <- as.logical(rbinom(10,1,0.5))
x2 <- sample(letters, 10, replace=TRUE)
x3 <- rnorm(10)
x4 <- ordered(cut(x3, -4:4, include.lowest=TRUE))
xx <- data.frame(x1, x2, x3, x4, stringsAsFactors = FALSE)

# matrix of distances between observations in xx
dx <- gower.dist(xx)
head(dx)

# matrix of distances between first obs. in xx
# and the remaining ones
gower.dist(data.x=xx[1:6,], data.y=xx[7:10,], var.weights = c(1,2,5,2))
```

| | |
|-------------|---|
| harmonize.x | <i>Harmonizes the marginal (joint) distribution of a set of variables observed independently in two sample surveys referred to the same target population</i> |
|-------------|---|

Description

This function permits to harmonize the marginal or the joint distribution of a set of variables observed independently in two sample surveys carried out on the same target population. This harmonization is carried out by using the calibration of the survey weights of the sample units in both the surveys according to the procedure suggested by Renssen (1998).

Usage

```
harmonize.x(svy.A, svy.B, form.x, x.tot=NULL,
           cal.method="linear", ...)
```

Arguments

| | |
|--------|--|
| svy.A | A svydesign R object that stores the data collected in the the survey A and all the information concerning the corresponding sampling design. This object can be created by using the function svydesign in the package survey . |
| svy.B | A svydesign R object that stores the data collected in the the survey B and all the information concerning the corresponding sampling design. This object can be created by using the function svydesign in the package survey . |
| form.x | A R formula specifying which of the variables, common to both the surveys, have to be considered, and how have to be considered. For instance <code>form.x=~x1+x2</code> means that the marginal distribution of the variables x1 and x2 have to be harmonized and there is also an 'Intercept'. In order to skip the intercept the formula has to be written in the following manner <code>form.x=~x1+x2-1</code> . |

When dealing with categorical variables, the formula `form.x~x1:x2-1` means that the harmonization has to be carried out by considering the joint distribution of the two variables (x_1 vs. x_2). To better understand how `form.x` works see [model.matrix](#) (see also [formula](#)).

Due to weights calibration features, it is preferable to work with categorical **X** variables. In some cases, the procedure may be successful when a single continuous variable is considered jointly with one or more categorical variables. When dealing with several continuous variable it may be preferable to categorize them.

- `x.tot` A vector or table with known population totals for the **X** variables. A vector is required when `cal.method="linear"` or `cal.method="raking"`. The names and the length of the vector depends on the way it is specified the argument `form.x` (see [model.matrix](#)). A contingency table is required when `cal.method="poststratify"` (for details see [postStratify](#)).
- When `x.tot` is not provided (i.e. `x.tot=NULL`) then the vector of totals is estimated as a weighted average of the totals estimated on the two surveys. The weight assigned to the totals estimated from A is $\lambda = n_A / (n_A + n_B)$; $1 - \lambda$ is the weight assigned to **X** totals estimated from survey B (n_A and n_B are the number of units in A and B respectively).
- `cal.method` A string that specifies how the calibration of the weights in `svy.A` and `svy.B` has to be carried out. By default linear calibration is performed (`cal.method="linear"`). In particular, the calibration is carried out by mean of the function [calibrate](#) in the package **survey**.
- Alternatively, it is possible to rake the origin survey weights by specifying `cal.method="raking"`. Finally, it is possible to perform a simple post-stratification by setting `cal.method="poststratify"`. Note that in this case the weights adjustments are carried out by considering the function [postStratify](#) in the package **survey**.
- ... Further arguments that may be necessary for calibration or post-stratification. The number of iterations used in calibration can be modified too by using the argument `maxit` (by default `maxit=50`). See [calibrate](#) for further details.

Details

This function harmonizes the totals of the **X** variables, observed in both survey A and survey B, to be equal to given known totals specified via `x.tot`. When these totals are not known (`x.tot=NULL`) they are estimated by combining the estimates derived from the two separate surveys. The harmonization is carried out according to a procedure suggested by Renssen (1998) based on calibration of survey weights (for major details on calibration see Sarndal and Lundstrom, 2005). The procedure is particularly suited to deal with categorical **X** variables, in this case it permits to harmonize the joint or the marginal distribution of the categorical variables being considered. Note that an incomplete crossing of the **X** variables can be considered: i.e. harmonisation wrt to the joint distribution of $X_1 \times X_2$ and the marginal distribution of X_3 .

The calibration procedure may not produce the final result due to convergence problems. In this case an error message appears. In order to reach convergence it may be necessary to launch the procedure

with less constraints (i.e a reduced number of population totals) by joining adjacent categories or by discarding some variables.

In some limited cases, it could be possible to consider both categorical and continuous variables. In this situation it may happen that calibration is not successful. In order to reach convergence it may be necessary to categorize the continuous **X** variables.

Post-stratification is a special case of calibration; all the weights of the units in a given post-stratum are modified so as to reproduce the known total for that post-stratum. Post-stratification avoids problems of convergence but, on the other hand, it may produce final weights with a higher variability than those derived from the calibration.

Value

A **R** with list the results of calibration procedures carried out on survey A and survey B, respectively. In particular the following components will be provided:

| | |
|-----------|---|
| cal.A | The survey object <code>svy.A</code> after the calibration; in particular, the weights now are calibrated with respect to the totals of the X variables. |
| cal.B | The survey object <code>svy.B</code> after the calibration; in particular, the weights now are calibrated with respect to the totals of the X variables. |
| weights.A | The new calibrated weights associated to the the units in <code>svy.A</code> . |
| weights.B | The new calibrated weights associated to the the units in <code>svy.B</code> . |
| call | Stores the call to this function with all the values specified for the various arguments (<code>call=match.call()</code>). |

Author(s)

Marcello D’Orazio <mdo.statmatch@gmail.com>

References

- D’Orazio, M., Di Zio, M. and Scanu, M. (2006). *Statistical Matching: Theory and Practice*. Wiley, Chichester.
- Renssen, R.H. (1998) “Use of Statistical Matching Techniques in Calibration Estimation”. *Survey Methodology*, N. **24**, pp. 171–183.
- Sarndal, C.E. and Lundstrom, S. (2005) *Estimation in Surveys with Nonresponse*. Wiley, Chichester.

See Also

[comb.samples](#), [calibrate](#), [svydesign](#), [postStratify](#),

Examples

```
data(quine, package="MASS") #loads quine from MASS
str(quine)

# split quine in two subsets
suppressWarnings(RNGversion("3.5.0"))
set.seed(7654)
```

```

lab.A <- sample(nrow(quine), 70, replace=TRUE)
quine.A <- quine[lab.A, c("Eth","Sex","Age","Lrn")]
quine.B <- quine[-lab.A, c("Eth","Sex","Age","Days")]

# create svydesign objects
require(survey)
quine.A$f <- 70/nrow(quine) # sampling fraction
quine.B$f <- (nrow(quine)-70)/nrow(quine)
svy.qA <- svydesign(~1, fpc=~f, data=quine.A)
svy.qB <- svydesign(~1, fpc=~f, data=quine.B)

#-----
# example (1)
# Harmonization of the distr. of Sex vs. Age
# usign poststratification

# (1.a) known population totals
# the population total are computed on the full data frame
tot.sex.age <- xtabs(~Sex+Age, data=quine)
tot.sex.age

out.hz <- harmonize.x(svy.A=svy.qA, svy.B=svy.qB, form.x=~Sex+Age,
  x.tot=tot.sex.age, cal.method="poststratify")

tot.A <- xtabs(out.hz$weights.A~Sex+Age, data=quine.A)
tot.B <- xtabs(out.hz$weights.B~Sex+Age, data=quine.B)

tot.sex.age-tot.A
tot.sex.age-tot.B

# (1.b) unknown population totals (x.tot=NULL)
# the population total is estimated by combining totals from the
# two surveys

out.hz <- harmonize.x(svy.A=svy.qA, svy.B=svy.qB, form.x=~Sex+Age,
  x.tot=NULL, cal.method="poststratify")

tot.A <- xtabs(out.hz$weights.A~Sex+Age, data=quine.A)
tot.B <- xtabs(out.hz$weights.B~Sex+Age, data=quine.B)

tot.A
tot.A-tot.B

#-----
# example (2)
# Harmonization wrt the maginal distribution
# of 'Eth', 'Sex' and 'Age'
# using linear calibration

# (2.a) vector of population total known
# estimated from the full data set
# note the formula! only marginal distribution of the
# variables are considered

```

```

tot.m <- colSums(model.matrix(~Eth+Sex+Age-1, data=quine))
tot.m

out.hz <- harmonize.x(svy.A=svy.qA, svy.B=svy.qB, x.tot=tot.m,
  form.x=~Eth+Sex+Age-1, cal.method="linear")

summary(out.hz$weights.A) #check for negative weights
summary(out.hz$weights.B) #check for negative weights

tot.m
svytable(formula=~Eth, design=out.hz$cal.A)
svytable(formula=~Eth, design=out.hz$cal.B)

svytable(formula=~Sex, design=out.hz$cal.A)
svytable(formula=~Sex, design=out.hz$cal.B)

# Note: margins are equal but joint distributions are not!
svytable(formula=~Sex+Age, design=out.hz$cal.A)
svytable(formula=~Sex+Age, design=out.hz$cal.B)

# (2.b) vector of population total unknown
out.hz <- harmonize.x(svy.A=svy.qA, svy.B=svy.qB, x.tot=NULL,
  form.x=~Eth+Sex+Age-1, cal.method="linear")
svytable(formula=~Eth, design=out.hz$cal.A)
svytable(formula=~Eth, design=out.hz$cal.B)

svytable(formula=~Sex, design=out.hz$cal.A)
svytable(formula=~Sex, design=out.hz$cal.B)

#-----
# example (3)
# Harmonization wrt the joint distribution of 'Sex' vs. 'Age'
# and the marginal distribution of 'Eth'
# using raking

# vector of population total known
# estimated from the full data set
# note the formula!
tot.m <- colSums(model.matrix(~Eth+(Sex:Age)-1, data=quine))
tot.m

out.hz <- harmonize.x(svy.A=svy.qA, svy.B=svy.qB, x.tot=tot.m,
  form.x=~Eth+(Sex:Age)-1, cal.method="raking")

summary(out.hz$weights.A) #check for negative weights
summary(out.hz$weights.B) #check for negative weights

tot.m
svytable(formula=~Eth, design=out.hz$cal.A)
svytable(formula=~Eth, design=out.hz$cal.B)

svytable(formula=~Sex+Age, design=out.hz$cal.A)
svytable(formula=~Sex+Age, design=out.hz$cal.B)

```

| | |
|------------------|--|
| mahalanobis.dist | <i>Computes the Mahalanobis Distance</i> |
|------------------|--|

Description

This function computes the Mahalanobis distance among units in a dataset or between observations in two distinct datasets.

Usage

```
mahalanobis.dist(data.x, data.y = NULL, vc = NULL, rob.vc = FALSE)
```

Arguments

| | |
|--------|--|
| data.x | <p>A matrix or a data frame containing variables that should be used in the computation of the distance between units. Only continuous variables are allowed. Missing values (NA) are not allowed.</p> <p>When only data.x is supplied, the distances between rows of data.x is computed.</p> |
| data.y | <p>A numeric matrix or data frame with the same variables, of the same type, as those in data.x (only continuous variables are allowed). Dissimilarities between rows of data.x and rows of data.y will be computed. If not provided, by default it is assumed data.y=data.x and only dissimilarities between rows of data.x will be computed.</p> |
| vc | <p>Covariance matrix that should be used in distance computation. If it is not supplied (vc = NULL) it is estimated from the input data. In particular, when vc = NULL and only data.x is supplied then the covariance matrix is estimated from data.x (i.e. vc = var(data.x)). On the contrary when vc = NULL and both data.x and data.y are available then the covariance matrix is estimated on the joined data sets.</p> |
| rob.vc | <p>Logical, if TRUE when vc = NULL the covariance matrix is estimated using a robust method by means of the function <code>cov.rob</code> in the package MASS.</p> |

Details

The Mahalanobis distance is calculated by means of:

$$d(i, j) = \sqrt{(x_i - x_j)^T S^{-1} (x_i - x_j)}$$

The covariance matrix S is estimated from the available data when `vc=NULL`, otherwise the one supplied via the argument `vc` is used.

Value

A matrix object with distances among rows of `data.x` and those of `data.y`.

Author(s)

Marcello D’Orazio <mdo.statmatch@gmail.com>

References

Mahalanobis, P C (1936) “On the generalised distance in statistics”. Proceedings of the National Institute of Sciences of India 2, pp. 49-55.

See Also

[mahalanobis](#)

Examples

```
md1 <- mahalanobis.dist(iris[1:6,1:4])
md2 <- mahalanobis.dist(data.x=iris[1:6,1:4], data.y=iris[51:60, 1:4])

vv <- var(iris[,1:4])
md1a <- mahalanobis.dist(data.x=iris[1:6,1:4], vc=vv)
md2a <- mahalanobis.dist(data.x=iris[1:6,1:4], data.y=iris[51:60, 1:4], vc=vv)
```

maximum.dist

Computes the Maximum Distance

Description

This function computes the Maximum distance (or L^∞ norm) between units in a dataset or between observations in two distinct datasets.

Usage

```
maximum.dist(data.x, data.y=data.x, rank=FALSE)
```

Arguments

`data.x` A matrix or a data frame containing variables that should be used in the computation of the distance. Only continuous variables are allowed. Missing values (NA) are not allowed.

When only `data.x` is supplied, the distances between rows of `data.x` are computed.

| | |
|--------|--|
| data.y | A numeric matrix or data frame with the same variables, of the same type, as those in data.x (only continuous variables are allowed). Dissimilarities between rows of data.x and rows of data.y will be computed. If not provided, by default it is assumed data.y=data.x and only dissimilarities between rows of data.x will be computed. |
| rank | Logical, when TRUE the original values are substituted by their ranks divided by the number of values plus one (following suggestion in Kovar et al. 1988). This rank transformation permits to remove the effect of different scales on the distance computation. When computing ranks the tied observations assume the average of their position (ties.method = "average" in calling the rank function). |

Details

This function computes the L^∞ distance also know as *minimax* distance. In practice the distance between two records is the maximum of the absolute differences on the available variables:

$$d(i, j) = \max(|x_{1i} - x_{1j}|, |x_{2i} - x_{2j}|, \dots, |x_{Ki} - x_{Kj}|)$$

When rank=TRUE the original values are substituted by their ranks divided by the number of values plus one (following suggestion in Kovar et al. 1988).

Value

A matrix object with distances between rows of data.x and those of data.y.

Author(s)

Marcello D’Orazio <mdo.statmatch@gmail.com>

References

Kovar, J.G., MacMillan, J. and Whitridge, P. (1988). “Overview and strategy for the Generalized Edit and Imputation System”. Statistics Canada, Methodology Branch Working Paper No. BSMD 88-007 E/F.

See Also

[rank](#),

Examples

```
md1 <- maximum.dist(iris[1:10,1:4])
md2 <- maximum.dist(iris[1:10,1:4], rank=TRUE)

md3 <- maximum.dist(data.x=iris[1:50,1:4], data.y=iris[51:100,1:4])
md4 <- maximum.dist(data.x=iris[1:50,1:4], data.y=iris[51:100,1:4], rank=TRUE)
```

mixed.mtc

*Statistical Matching via Mixed Methods***Description**

This function implements some mixed methods to perform statistical matching between two data sources.

Usage

```
mixed.mtc(data.rec, data.don, match.vars, y.rec, z.don, method="ML",
          rho.yz=NULL, micro=FALSE, constr.alg="Hungarian")
```

Arguments

| | |
|------------|--|
| data.rec | A matrix or data frame that plays the role of <i>recipient</i> in the statistical matching application. This data set must contain all variables (columns) that should be used in statistical matching, i.e. the variables called by the arguments <code>match.vars</code> and <code>y.rec</code> . Note that continuous variables are expected, if there are some categorical variables they are re-coded into dummies. Missing values (NA) are not allowed. |
| data.don | A matrix or data frame that plays the role of <i>donor</i> in the statistical matching application. This data set must contain all the numeric variables (columns) that should be used in statistical matching, i.e. the variables called by the arguments <code>match.vars</code> and <code>z.don</code> . Note that continuous variables are expected, if there are some categorical variables they are re-coded into dummies. Missing values (NA) are not allowed. |
| match.vars | A character vector with the names of the common variables (the columns in both the data frames) to be used as matching variables (X). |
| y.rec | A character vector with the name of the target variable Y that is observed only for units in <code>data.rec</code> . Only one continuous variable is allowed. |
| z.don | A character vector with the name of the target variable Z that is observed only for units in <code>data.don</code> . Only one continuous variable is allowed. |
| method | A character vector that identifies the method that should be used to estimate the parameters of the regression models: Y vs. X and Z vs. X . Maximum Likelihood method is used when <code>method="ML"</code> (default); on the contrary, when <code>method="MS"</code> the parameters are estimated according to approach proposed by Moriarity and Scheuren (2001 and 2003). See Details for further information. |
| rho.yz | A numeric value representing a guess for the correlation between the Y (<code>y.rec</code>) and the Z variable (<code>z.don</code>) that are not jointly observed. When <code>method="MS"</code> then the argument <code>cor.yz</code> must specify the value of the correlation coefficient ρ_{YZ} ; on the contrary, when <code>method="ML"</code> , it must specify the <i>partial correlation coefficient</i> between Y and Z given X ($\rho_{YZ X}$). |

| | |
|-------------------------|---|
| | By default (<code>rho.yz=NULL</code>). In practice, in absence of auxiliary information concerning the correlation coefficient or the partial correlation coefficient, the statistical matching is carried out under the assumption of independence between Y and Z given \mathbf{X} (Conditional Independence Assumption, CIA), i.e. $\rho_{YZ \mathbf{X}} = 0$. |
| <code>micro</code> | Logical. When <code>micro=FALSE</code> (default) only the parameters' estimates are returned. On the contrary, when <code>micro=TRUE</code> the function returns also <code>data.rec</code> filled in with the values for the variable Z . The donors for filling in Z in <code>data.rec</code> are identified using a constrained distance hot deck method. In this case, the number of units (rows) in <code>data.don</code> must be greater or equal to the number of units (rows) in <code>data.rec</code> . See next argument and Details for further information. |
| <code>constr.alg</code> | A string that has to be specified when <code>micro=TRUE</code> , in order to solve the transportation problem involved by the constrained distance hot deck method. Two choices are available: "lpSolve" and "Hungarian". In the first case, <code>constr.alg="lpSolve"</code> , the transportation problem is solved by means of the function <code>lp.transport</code> available in the package <code>lpSolve</code> . When <code>constr.alg="Hungarian"</code> (default) the transportation problem is solved using the Hungarian method implemented in the function <code>solve_LSAP</code> available in the package <code>clue</code> (Hornik, 2012). Note that Hungarian algorithm is more efficient and requires less processing time. |

Details

This function implements some mixed methods to perform statistical matching. A mixed method consists of two steps:

- (1) adoption of a parametric model for the joint distribution of (\mathbf{X}, Y, Z) and estimation of its parameters;
- (2) derivation of a complete "synthetic" data set (recipient data set filled in with values for the Z variable) using a nonparametric approach.

In this case, as far as (1) is concerned, it is assumed that (\mathbf{X}, Y, Z) follows a multivariate normal distribution. Please note that if some of the \mathbf{X} are categorical, then they are recoded into dummies before starting with the estimation. In such a case, the assumption of multivariate normal distribution may be questionable.

The whole procedure is based on the imputation method known as *predictive mean matching*. The procedure consists of three steps:

step 1a) Regression step: the two linear regression models Y vs. \mathbf{X} and Z vs. \mathbf{X} are considered and their parameters are estimated.

step 1b) Computation of intermediate values. For the units in `data.rec` the following intermediate values are derived:

$$\tilde{z}_a = \hat{\alpha}_Z + \hat{\beta}_{Z\mathbf{X}}\mathbf{x}_a + e_a$$

for each $a = 1, \dots, n_A$, being n_A the number of units in `data.rec` (rows of `data.rec`). Note that, e_a is a random draw from the multivariate normal distribution with zero mean and estimated residual variance $\hat{\sigma}_{Z|\mathbf{X}}$.

Similarly, for the units in `data.don` the following intermediate values are derived:

$$\tilde{y}_b = \hat{\alpha}_Y + \hat{\beta}_{Y\mathbf{X}}\mathbf{x}_b + e_b$$

for each $b = 1, \dots, n_B$, being n_B the number of units in data.don (rows of data.don). e_b is a random draw from the multivariate normal distribution with zero mean and estimated residual variance $\hat{\sigma}_{Y|\mathbf{X}}$.

step 2) Matching step. For each observation (row) in data.rec a donor is chosen in data.don through a nearest neighbor constrained distance hot deck procedure. The distances are computed between (y_a, \tilde{z}_a) and (\tilde{y}_b, z_b) using Mahalanobis distance.

For further details see Sections 2.5.1 and 3.6.1 in D’Orazio *et al.* (2006).

In step 1a) the parameters of the regression model can be estimated by means of the Maximum Likelihood method (method="ML") (see D’Orazio *et al.*, 2006, pp. 19–23,73–75) or, using the Moriarity and Scheuren (2001 and 2003) approach (method="MS") (see also D’Orazio *et al.*, 2006, pp. 75–76). The two estimation methods are compared in D’Orazio *et al.* (2005).

When method="MS", if the value specified for the argument rho.yz is not compatible with the other correlation coefficients estimated from the data, then it is substituted with the closest value compatible with the other estimated coefficients.

When micro=FALSE only the estimation of the parameters is performed (step 1a). Otherwise, (micro=TRUE) the whole procedure is carried out.

Value

A list with a varying number of components depending on the values of the arguments method and rho.yz.

| | |
|---------------|---|
| mu | The estimated mean vector. |
| vc | The estimated variance–covariance matrix. |
| cor | The estimated correlation matrix. |
| res.var | A vector with estimates of the residual variances $\sigma_{Y Z\mathbf{X}}$ and $\sigma_{Z Y\mathbf{X}}$. |
| start.prho.yz | It is the initial guess for the partial correlation coefficient $\rho_{YZ \mathbf{X}}$ passed in input via the rho.yz argument when method="ML". |
| rho.yz | Returned in output only when method="MS". It is a vector with four values: the initial guess for ρ_{YZ} ; the lower and upper bounds for $\hat{\rho}_{YZ}$ in the statistical matching framework given the correlation coefficients between Y and X and the correlation coefficients between Z and X estimated from the available data; and, finally, the closest admissible value used in computations instead of the initial rho.yz that resulted not coherent with the others correlation coefficients estimated from the available data. |
| phi | When method="MS". Estimates of the ϕ terms introduced by Moriarity and Scheuren (2001 and 2003). |
| filled.rec | The data.rec filled in with the values of Z. It is returned only when micro=TRUE. |
| mtc.ids | when micro=TRUE. This is a matrix with the same number of rows of data.rec and two columns. The first column contains the row names of the data.rec and the second column contains the row names of the corresponding donors selected from the data.don. When the input matrices do not contain row names, a numeric matrix with the indexes of the rows is provided. |

dist.rd A vector with the distances between each recipient unit and the corresponding donor, returned only in case micro=TRUE.

call How the function has been called.

Author(s)

Marcello D’Orazio <mdo.statmatch@gmail.com>

References

D’Orazio, M., Di Zio, M. and Scanu, M. (2005). “A comparison among different estimators of regression parameters on statistically matched files through an extensive simulation study”, *Contributi*, **2005/10**, Istituto Nazionale di Statistica, Rome.

D’Orazio, M., Di Zio, M. and Scanu, M. (2006). *Statistical Matching: Theory and Practice*. Wiley, Chichester.

Hornik K. (2012). clue: Cluster ensembles. R package version 0.3-45. <https://CRAN.R-project.org/package=clue>.

Moriarity, C., and Scheuren, F. (2001). “Statistical matching: a paradigm for assessing the uncertainty in the procedure”. *Journal of Official Statistics*, **17**, 407–422.

Moriarity, C., and Scheuren, F. (2003). “A note on Rubin’s statistical matching using file concatenation with adjusted weights and multiple imputation”, *Journal of Business and Economic Statistics*, **21**, 65–73.

See Also

[NND.hotdeck](#), [mahalanobis.dist](#)

Examples

```
# reproduce the statistical matching framework
# starting from the iris data.frame
suppressWarnings(RNGversion("3.5.0"))
set.seed(98765)
pos <- sample(1:150, 50, replace=FALSE)
ir.A <- iris[pos,c(1,3:5)]
ir.B <- iris[-pos, 2:5]

xx <- intersect(colnames(ir.A), colnames(ir.B))
xx # common variables

# ML estimation method under CIA ((rho_YZ|X=0));
# only parameter estimates (micro=FALSE)
# only continuous matching variables
xx.mtc <- c("Petal.Length", "Petal.Width")
mtc.1 <- mixed.mtc(data.rec=ir.A, data.don=ir.B, match.vars=xx.mtc,
                  y.rec="Sepal.Length", z.don="Sepal.Width")

# estimated correlation matrix
mtc.1$cor
```

```

# ML estimation method under CIA ((rho_YZ|X=0));
# only parameter estimates (micro=FALSE)
# categorical variable 'Species' used as matching variable

xx.mtc <- xx
mtc.2 <- mixed.mtc(data.rec=ir.A, data.don=ir.B, match.vars=xx.mtc,
                  y.rec="Sepal.Length", z.don="Sepal.Width")

# estimated correlation matrix
mtc.2$cor

# ML estimation method with partial correlation coefficient
# set equal to 0.5 (rho_YZ|X=0.5)
# only parameter estimates (micro=FALSE)

mtc.3 <- mixed.mtc(data.rec=ir.A, data.don=ir.B, match.vars=xx.mtc,
                  y.rec="Sepal.Length", z.don="Sepal.Width",
                  rho.yz=0.5)

# estimated correlation matrix
mtc.3$cor

# ML estimation method with partial correlation coefficient
# set equal to 0.5 (rho_YZ|X=0.5)
# with imputation step (micro=TRUE)

mtc.4 <- mixed.mtc(data.rec=ir.A, data.don=ir.B, match.vars=xx.mtc,
                  y.rec="Sepal.Length", z.don="Sepal.Width",
                  rho.yz=0.5, micro=TRUE, constr.alg="Hungarian")

# first rows of data.rec filled in with z
head(mtc.4$filled.rec)

#
# Moriarity and Scheuren estimation method under CIA;
# only with parameter estimates (micro=FALSE)
mtc.5 <- mixed.mtc(data.rec=ir.A, data.don=ir.B, match.vars=xx.mtc,
                  y.rec="Sepal.Length", z.don="Sepal.Width",
                  method="MS")

# the starting value of rho.yz and the value used
# in computations
mtc.5$rho.yz

# estimated correlation matrix
mtc.5$cor

# Moriarity and Scheuren estimation method
# with correlation coefficient set equal to -0.15 (rho_YZ=-0.15)
# with imputation step (micro=TRUE)

mtc.6 <- mixed.mtc(data.rec=ir.A, data.don=ir.B, match.vars=xx.mtc,

```

```

y.rec="Sepal.Length", z.don="Sepal.Width",
method="MS", rho.yz=-0.15,
micro=TRUE, constr.alg="lpSolve")

# the starting value of rho.yz and the value used
# in computations
mtc.6$rho.yz

# estimated correlation matrix
mtc.6$cor

# first rows of data.rec filled in with z imputed values
head(mtc.6$filled.rec)

```

NND.hotdeck

Distance Hot Deck method.

Description

This function implements the distance hot deck method to match the records of two data sources that share some variables.

Usage

```

NND.hotdeck(data.rec, data.don, match.vars,
             don.class=NULL, dist.fun="Manhattan",
             constrained=FALSE, constr.alg="Hungarian",
             k=1, keep.t=FALSE, ...)

```

Arguments

| | |
|------------|---|
| data.rec | A matrix or data frame that plays the role of <i>recipient</i> . This data frame must contain the variables (columns) that should be used, directly or indirectly, in the matching application (specified via <code>match.vars</code> and eventually <code>don.class</code>). Missing values (NA) are allowed. |
| data.don | A matrix or data frame that plays the role of <i>donor</i> . The variables (columns) involved, directly or indirectly, in the computation of distance must be the same and of the same type as those in <code>data.rec</code> (specified via <code>match.vars</code> and eventually <code>don.class</code>). |
| match.vars | A character vector with the names of the matching variables (the columns in both the data frames) that have to be used to compute distances between records (rows) in <code>data.rec</code> and those in <code>data.don</code> . The variables used in computing distances may contain missing values but only a limited number of distance functions can handle them (see Details for clarifications). |

| | |
|--------------------------|---|
| <code>don.class</code> | <p>A character vector with the names of the variables (columns in both the data frames) that have to be used to identify the donation classes. In this case the computation of distances is limited to those units of <code>data.rec</code> and <code>data.doc</code> that belong to the same donation class. The case of empty donation classes should be avoided. It would be preferable that variables used to form donation classes are defined as factor.</p> <p>The variables chosen for the creation of the donation classes should NOT contain missing values (NAs).</p> <p>When not specified (default), no donation classes are used. This choice may require more memory to store a larger distance matrix and a higher computational effort.</p> |
| <code>dist.fun</code> | <p>A string with the name of the distance function that has to be used. The following distances are allowed: “Manhattan” (aka “City block”; default), “Euclidean”, “Mahalanobis”, “exact” or “exact matching”, “Gower”, “minimax” or one of the distance functions available in the package proxy. Note that the distance is computed using the function <code>dist</code> of the package proxy with the exception of the “Gower” (see function <code>gower.dist</code> for details), “Mahalanobis” (function <code>mahalanobis.dist</code>) and “minimax” (see <code>maximum.dist</code>) cases.</p> <p>When <code>dist.fun="Manhattan"</code>, “Euclidean”, “Mahalanobis” or “minimax” all the matching variables in <code>data.rec</code> and <code>data.doc</code> must be numeric. When <code>dist.fun="exact"</code> or <code>dist.fun="exact matching"</code>, all the variables in <code>data.rec</code> and <code>data.doc</code> will be converted to character and, as far as the distance computation is concerned, they will be treated as categorical nominal variables, i.e. distance is 0 if a couple of units presents the same response category and 1 otherwise.</p> |
| <code>constrained</code> | <p>Logical. When <code>constrained=FALSE</code> (default) each record in <code>data.doc</code> can be used as a donor more than once. On the contrary, when <code>constrained=TRUE</code> each record in <code>data.doc</code> can be used as a donor only <code>k</code> times. In this case, the set of donors is selected by solving an optimization problem, whose goal is to minimize the overall matching distance. See description of the argument <code>constr.alg</code> for details.</p> |
| <code>constr.alg</code> | <p>A string that has to be specified when <code>constrained=TRUE</code>. Two choices are available: “lpSolve” and “hungarian”. In the first case, <code>constr.alg="lpSolve"</code>, the optimization problem is solved by means of the function <code>lp.transport</code> available in the package lpSolve. When <code>constr.alg="hungarian"</code> (default) the problem is solved using the Hungarian method, implemented in function <code>solve_LSAP</code> available in the package clue. Note that Hungarian algorithm is faster and more efficient if compared to <code>constr.alg="lpSolve"</code> but it allows selecting a donor just once, i.e. <code>k = 1</code>.</p> |
| <code>k</code> | <p>The number of times that a unit in <code>data.doc</code> can be selected as a donor when <code>constrained=TRUE</code> (default <code>k = 1</code>). When <code>k>1</code> then optimization problem can be solved by setting <code>constr.alg="lpSolve"</code>. Hungarian algorithm (<code>constr.alg="hungarian"</code>) can be used only when <code>k = 1</code>.</p> |
| <code>keep.t</code> | <p>Logical, when donation classes are used by setting <code>keep.t=TRUE</code> prints information on the donation classes being processed (by default <code>keep.t=FALSE</code>).</p> |
| <code>...</code> | <p>Additional arguments that may be required by <code>gower.dist</code>, <code>maximum.dist</code>, or <code>dist</code>.</p> |

Details

This function finds a donor record in `data.don` for each record in `data.rec`. In the unconstrained case, it searches for the closest donor record according to the chosen distance function. When for a given recipient record there are more donors available at the minimum distance, one of them is picked at random.

In the constrained case a donor can be used just a fixed number of times, as specified by the `k` argument, but the whole set of donors is chosen in order to minimize the overall matching distance. When `k=1` the number of units (rows) in the donor data set has to be larger or equal to the number of units of the recipient data set; when the donation classes are used, this condition must be satisfied in each donation class. For further details on nearest neighbor distance hot deck refer to Chapter 2 in D’Orazio *et al.* (2006).

This function can also be used to impute missing values in a data set using the nearest neighbor distance hot deck. In this case `data.rec` is the part of the initial data set that contains missing values on the target variable; on the contrary, `data.don` is the part of the data set without missing values on it. See R code in the Examples for details.

Please note that only “Gower” and “minimax” distance functions allow for the presence of missing values (NAs) in the variables used in computing distances. In both the cases when one of the of the observations presents a variable showing an NA, then this variable is excluded from the computation of distance between them.

Value

A R list with the following components:

| | |
|----------------------|--|
| <code>mtc.ids</code> | A matrix with the same number of rows of <code>data.rec</code> and two columns. The first column contains the row names of the <code>data.rec</code> and the second column contains the row names of the corresponding donors selected from the <code>data.don</code> . When the input matrices do not contain row names, a numeric matrix with the indexes of the rows is provided. |
| <code>dist.rd</code> | A vector with the distances between each recipient unit and the corresponding donor. |
| <code>noad</code> | When <code>constrained=FALSE</code> , it reports the number of available donors at the minimum distance for each recipient unit. |
| <code>call</code> | How the function has been called. |

Author(s)

Marcello D’Orazio <mdo.statmatch@gmail.com>

References

- D’Orazio, M., Di Zio, M. and Scanu, M. (2006). *Statistical Matching: Theory and Practice*. Wiley, Chichester.
- Hornik K. (2012). `clue`: Cluster ensembles. R package version 0.3-45. <https://CRAN.R-project.org/package=clue>.
- Rodgers, W.L. (1984). “An evaluation of statistical matching”. *Journal of Business and Economic Statistics*, **2**, 91–102.

Singh, A.C., Mantel, H., Kinack, M. and Rowe, G. (1993). "Statistical matching: use of auxiliary information as an alternative to the conditional independence assumption". *Survey Methodology*, **19**, 59–79.

See Also

[RANDwNND.hotdeck](#)

Examples

```
# create the classical matching framework
lab <- c(1:15, 51:65, 101:115)
iris.rec <- iris[lab, c(1:3,5)] # recipient data.frame
iris.don <- iris[-lab, c(1:2,4:5)] #donor data.frame

# Now iris.rec and iris.don have the variables
# "Sepal.Length", "Sepal.Width" and "Species"
# in common.
# "Petal.Length" is available only in iris.rec
# "Petal.Width" is available only in iris.don

# Find the closest donors computing distance
# on "Sepal.Length" and "Sepal.Width"
# unconstrained case, Euclidean distance

out.NND.1 <- NND.hotdeck(data.rec=iris.rec, data.don=iris.don,
                        match.vars=c("Sepal.Length", "Sepal.Width") )

# create the synthetic data.set:
# fill in "Petal.Width" in iris.rec

fused.1 <- create.fused(data.rec=iris.rec, data.don=iris.don,
                       mtc.ids=out.NND.1$mtc.ids, z.vars="Petal.Width")
head(fused.1)

# Find the closest donors computing distance
# on "Sepal.Length", "Sepal.Width" and Species;
# unconstrained case, Gower's distance

out.NND.2 <- NND.hotdeck(data.rec=iris.rec, data.don=iris.don,
                        match.vars=c("Sepal.Length", "Sepal.Width", "Species"),
                        dist.fun="Gower")

# find the closest donors using "Species" to form donation classes
# and "Sepal.Length" and "Sepal.Width" to compute distance;
# unconstrained case.

out.NND.3 <- NND.hotdeck(data.rec=iris.rec, data.don=iris.don,
                        match.vars=c("Sepal.Length", "Sepal.Width"),
                        don.class="Species")
```

```

# find the donors using "Species" to form donation classes
# and "Sepal.Length" and "Sepal.Width" to compute distance;
# constrained case, "Hungarian" algorithm

library(clue)
out.NND.4 <- NND.hotdeck(data.rec=iris.rec, data.don=iris.don,
                        match.vars=c("Sepal.Length", "Sepal.Width"),
                        don.class="Species", constrained=TRUE,
                        constr.alg="Hungarian")

# Example of Imputation of missing values.
# Introducing missing values in iris
ir.mat <- iris
miss <- rbinom(nrow(iris), 1, 0.3)
ir.mat[miss==1,"Sepal.Length"] <- NA
iris.rec <- ir.mat[miss==1,-1]
iris.don <- ir.mat[miss==0,]

#search for NND donors
imp.NND <- NND.hotdeck(data.rec=iris.rec, data.don=iris.don,
                      match.vars=c("Sepal.Width", "Petal.Length", "Petal.Width"),
                      don.class="Species")

# imputing missing values
iris.rec.imp <- create.fused(data.rec=iris.rec, data.don=iris.don,
                           mtc.ids=imp.NND$mtc.ids, z.vars="Sepal.Length")

# rebuild the imputed data.frame
final <- rbind(iris.rec.imp, iris.don)
head(final)

```

pBayes

Pseudo-Bayes estimates of cell probabilities

Description

Estimation of cells counts in contingency tables by means of the pseudo-Bayes estimator.

Usage

```
pBayes(x, method="m.ind", const=NULL)
```

Arguments

x A contingency table with observed cell counts. Typically the output of `table` or `xtabs`. More in general an R array with the counts.

| | |
|--------|--|
| method | <p>The method for estimating the final cell frequencies. The following options are available:</p> <p>method = "Jeffreys", consists in adding 0.5 to each cell before estimation of the relative frequencies.</p> <p>method = "minimax", consists in adding $\sqrt{(n)}/c$ to each cell before estimation of the relative frequencies, being n the sum of all the counts and c the number of cells in the table.</p> <p>method = "invcat", consists in adding $1/c$ to each cell before estimation of the relative frequencies.</p> <p>method = "user", consists in adding a user defined constant a ($a > 0$) to each cell before estimation of the relative frequencies. The constant a should be passed via the argument <code>const</code>.</p> <p>method = "m.ind", the prior guess for the unknown cell probabilities is obtained by considering estimated probabilities under the mutual independence hypothesis. This option is available when dealing with at least two-way contingency tables (<code>length(dim(x))>=2</code>).</p> <p>method = "h.assoc", the prior guess for the unknown cell probabilities is obtained by considering estimated probabilities under the homogeneous association hypothesis. This option is available when dealing with at least two-way contingency tables (<code>length(dim(x))>=2</code>).</p> |
| const | <p>Numeric value, a user defined constant a ($a > 0$) to be added to each cell before estimation of the relative frequencies when <code>method = "user"</code>. As a general rule of thumb, it is preferable to avoid that the sum of constant over all the cells is greater than $0.20 \times n$.</p> |

Details

This function estimates the frequencies in a contingency table by using the pseudo-Bayes approach. In practice the estimator being considered is a weighted average of the input (observed) cells counts n_h and a suitable prior guess, γ_h , for cells probabilities :

$$\tilde{p}_h = \frac{n}{n+K} \hat{p}_h + \frac{K}{n+K} \gamma_h$$

K depends on the parameters of Dirichlet prior distribution being considered (for major details see Chapter 12 in Bishop et al., 1974). It is worth noting that with a constant prior guess $\gamma_h = 1/c$ ($h = 1, 2, \dots, c$), then $K = 1$ and in practice corresponds to adding $1/c$ to each cell before estimation of the relative frequencies (`method = "invcat"`); $K = c/2$ when the constant 0.5 is added to each cell (`method = "Jeffreys"`); finally $K = \sqrt{n}$ when the quantity \sqrt{n}/c is added to each cell (`method = "minimax"`). All these cases corresponds to adding a flattening constant; the higher is the value of K the more the estimates will be shrunk towards $\gamma_h = 1/c$ (flattening).

When `method = "m.ind"` the prior guess γ_h is estimated under the hypothesis of mutual independence between the variables crossed in the initial contingency table x , supposed to be at least a two-way table. In this case the value of K is estimated via a data-driven approach by considering

$$\hat{K} = \frac{1 - \sum_h \hat{p}_h^2}{\sum_h (\hat{\gamma}_h - \hat{p}_h)^2}$$

On the contrary, when `method = "h.assoc"` the prior guess γ_h is estimated under the hypothesis of homogeneous association between the variables crossed in the initial contingency table `x`.

Please note that when the input table is estimated from sample data where a weight is assigned to each unit, the weights should be used in estimating the input table, but it is suggested to rescale them so that their sum is equal to n , the sample size.

Value

A list object with three components.

| | |
|----------------------|--|
| <code>info</code> | A vector with the sample size " <code>n</code> ", the number of cells (" <code>no.cells</code> ") in <code>x</code> , the average cell frequency (" <code>av.cfr</code> "), the number of cells showing frequencies equal to zero (" <code>no.0s</code> "), the <code>const</code> input argument, the chosen/estimated K (" <code>K</code> ") and the relative size of K , i.e. $K/(n + K)$ (" <code>rel.K</code> "). |
| <code>prior</code> | A table having the same dimension as <code>x</code> with the considered prior values for the cell frequencies. |
| <code>pseudoB</code> | A table with having the same dimension as <code>x</code> providing the pseudo-Bayes estimates for the cell frequencies in <code>x</code> . |

Author(s)

Marcello D'Orazio <mdo.statmatch@gmail.com>

References

Bishop Y.M.M., Fienberg, S.E., Holland, P.W. (1974) *Discrete Multivariate Analysis: Theory and Practice*. The Massachusetts Institute of Technology

Examples

```
data(samp.A, package="StatMatch")
tab <- xtabs(~ area5 + urb + c.age + sex + edu7, data = samp.A)
out.pb <- pBayes(x=tab, method="m.ind")
out.pb$info

out.pb <- pBayes(x=tab, method="h.assoc")
out.pb$info

out.pb <- pBayes(x=tab, method="Jeffreys")
out.pb$info

# usage of weights in estimating the input table
n <- nrow(samp.A)
r.w <- samp.A$ww / sum(samp.A$ww) * n # rescale weights to sum up to n
tab.w <- xtabs(r.w ~ area5 + urb + c.age + sex + edu7, data = samp.A)
out.pbw <- pBayes(x=tab.w, method="m.ind")
out.pbw$info
```

| | |
|------------|---|
| plotBounds | <i>Graphical representation of the uncertainty bounds estimated through the Frechet.bounds.cat function</i> |
|------------|---|

Description

The function uses the output of the function [Frechet.bounds.cat](#) to produce a basic graphical representation of the uncertainty bounds related to the contingency table of Y vs. Z.

Usage

```
plotBounds(outFB)
```

Arguments

outFB the list provided in output from [Frechet.bounds.cat](#).

Details

This function represents graphically the uncertainty bounds estimated by the function [Frechet.bounds.cat](#) for each relative frequency in the contingency table of Y vs. Z. the dotted line indicates the width of the bounds estimated without conditioning on the Xs (the size is reported in parenthesis below the line). The full line indicates the width of the estimated bounds conditional on the Xs (expected conditional Frechet bounds for the relative frequencies in the contingency table of Y vs. Z (size reported below the line, not in the the parenthesis). Not that when the X are not used it is drawn only the width of the unconditional bounds and the size is shown below the line.

The figure on the top od the line indicated the estimated relative frequency under the assumption of independence between Y and Z conditional on one or more X variables (Conditional Independence Assumption, CIA; for details see help pages of [Frechet.bounds.cat](#)), otherwise it corresponds to the estimated relative frequency under the assumption of independence between Y and Z.

Value

The required graphical representation is drawn using standard **graphics** facilities.

Author(s)

Marcello D’Orazio <mdo.statmatch@gmail.com>

See Also

[Frechet.bounds.cat](#)

Examples

```
# # compute the tables required by Frechet.bounds.cat()
# freq.xA <- xtabs(~sex+c.age, data=samp.A)
# freq.xB <- xtabs(~sex+c.age, data=samp.B)
# freq.xy <- xtabs(~sex+c.age+c.neti, data=samp.A)
# freq.xz <- xtabs(~sex+c.age+labour5, data=samp.B)
#
# # apply Frechet.bounds.cat()
# bounds.yz <- Frechet.bounds.cat(tab.x=freq.xA+freq.xB, tab.xy=freq.xy,
#                                tab.xz=freq.xz, print.f="data.frame")
#
#
# plot.bounds(bounds.yz)
```

| | |
|----------|--|
| plotCont | <i>graphical comparison of the estimated distributions for the same continuous variable.</i> |
|----------|--|

Description

Compares graphically the estimated distributions for the same continuous variable using data coming from two different data sources.

Usage

```
plotCont(data.A, data.B, xlab.A, xlab.B=NULL, w.A=NULL, w.B=NULL,
         type="density", ref=FALSE)
```

Arguments

| | |
|--------|--|
| data.A | A dataframe or matrix containing the variable of interest xlab.A and eventual associated survey weights w.A. |
| data.B | A dataframe or matrix containing the variable of interest xlab.B and eventual associated survey weights w.B. |
| xlab.A | Character string providing the name of the variable in data.A whose distribution should be represented graphically and compared with that estimated from data.B. |
| xlab.B | Character string providing the name of the variable in data.B whose distribution should be represented graphically and compared with that estimated from data.A. If xlab.B=NULL (default) then it assumed xlab.B=xlab.A. |
| w.A | Character string providing the name of the optional weighting variable in data.A that, in case, should be used to estimate the distribution of xlab.A |
| w.B | Character string providing the name of the optional weighting variable in data.B that, in case, should be used to estimate the distribution of xlab.B |

| | |
|------|--|
| type | A character string indicating the type of graphical representation that should be used to compare the estimated distributions of <code>xlab.A</code> and <code>xlab.B</code> . By default (<code>type="density"</code>) density plots are used. Other possible options are “ <code>ecdf</code> ”, “ <code>qqplot</code> ”, “ <code>qqshift</code> ” and “ <code>hist</code> ”. See Details for more information. |
| ref | Logical, indicating whether the distribution estimated from <code>data.B</code> should be considered the reference or not. Default <code>ref=FALSE</code> . when Default <code>ref=TRUE</code> the estimation of the histograms, the density and the empirical cumulative distribution function are guided by data in <code>data.B</code> |

Details

This function graphically compares the distribution of the same variable, but estimated from data from two different data sources. The graphical comparison can be performed in several ways. With `type="hist"`, the continuous variable is categorized and the corresponding histograms estimated from `data.A` and `data.B` are compared. If present, the weights are used to estimate the relative frequencies. Note that the breaks to categorize the variable are decided according to the Freedman-Diaconis rule (`nclass`), and in this case with `ref=TRUE` the IQR is estimated from `data.B` alone, while with `ref=FALSE` it is estimated by joining the two data sources.

With `type="density"` the density plots are drawn; when available the weights are used in the estimation of the density based on the histograms (as suggested by Bellhouse and Stafford, 1999). When `type="ecdf"` the comparison relies on the empirical cumulative distribution function, that can be estimated considering the weights. Note that when `ref=TRUE` the estimation of the density and the empirical cumulative distribution are guided by the data in `data.B`.

The comparison is based on percentiles with `type="qqplot"` and `type="qqshift"`. In the first case, the function draws a scatterplot (red dots) of the estimated percentiles of `xlab.A` against those of `xlab.B`; the dashed line indicates the ideal situation of equality of percentiles (points lying on the line). When `type="qqshift"`, the scatterplot refers to $(\text{percentiles.A} - \text{percentiles.B})$ vs. `percentiles.B`; in this case, the points lying on the horizontal line passing through 0 indicate equality (difference equal to 0). Note that the number of estimated percentiles depends on the minimum between the two sample sizes. Only quartiles are calculated if $\min(n.A, n.B) \leq 50$; quintiles are estimated if $\min(n.A, n.B) > 50$ and $\min(n.A, n.B) \leq 150$; deciles are estimated if $\min(n.A, n.B) > 150$ and $\min(n.A, n.B) \leq 250$; finally, quantiles for `probs=seq(from = 0.05, to = 0.95, by = 0.05)` are estimated when $\min(n.A, n.B) > 250$. If survey weights are available (indicated by `w.A` and/or `w.B`), they are used to estimate the quantiles by calling the function `wtd.quantile` in the package **Hmisc**.

Value

The required graphical representation is drawn using the **ggplot2** facilities.

Author(s)

Marcello D’Orazio <mdo.statmatch@gmail.com>

References

Bellhouse D.R. and J. E. Stafford (1999). “Density Estimation from Complex Surveys”. *Statistica Sinica*, **9**, 407–424.

See Also[comp.cont](#)**Examples**

```
# plotCont(data.A = samp.A, data.B = samp.B, xlab.A="age")
# plotCont(data.A = samp.A, data.B = samp.B, xlab.A="age", w.A = "ww")
```

| | |
|---------|---|
| plotTab | <i>Graphical comparison of the estimated distributions for the same categorical variable.</i> |
|---------|---|

Description

Compares graphically the estimated distributions for the same categorical variable(s) using data coming from two different data sources.

Usage

```
plotTab(data.A, data.B, xlab.A, xlab.B=NULL, w.A=NULL, w.B=NULL)
```

Arguments

| | |
|--------|--|
| data.A | A dataframe or matrix containing the variable of interest xlab.A and eventual associated survey weights w.A. |
| data.B | A dataframe or matrix containing the variable of interest xlab.B and eventual associated survey weights w.B. |
| xlab.A | Character string providing the name(s) of one or more variables in data.A whose (joint) distribution should be represented graphically and compared with that estimated from data.B. |
| xlab.B | Character string providing the name(s) of one or more variables in data.A whose (joint) distribution should be represented graphically and compared with that estimated from data.A. If xlab.B=NULL (default) then it assumed xlab.B=xlab.A. |
| w.A | Character string providing the name of the optional weighting variable in data.A that, in case, should be used to estimate the distribution of xlab.A |
| w.B | Character string providing the name of the optional weighting variable in data.B that, in case, should be used to estimate the distribution of xlab.B |

Details

This function compares graphically the (joint) distribution of the same variables but estimated from data coming from two different data sources. The graphical comparison is done using barcharts.

Value

The required graphical representation is drawn using the **ggplot2** facilities.

Author(s)

Marcello D’Orazio <mdo.statmatch@gmail.com>

See Also

[comp.prop](#)

Examples

```
# plotTab(data.A = samp.A, data.B = samp.B, xlab.A="edu7", w.A = "ww")
# plotTab(data.A = samp.A, data.B = samp.B, xlab.A=c("urb", "sex"), w.A = "ww", w.B="ww")
```

pw.assoc

Pairwise measures between categorical variables

Description

This function computes some association and *Proportional Reduction in Error* (PRE) measures between a categorical nominal variable and each of the other available predictors (being also categorical variables).

Usage

```
pw.assoc(formula, data, weights=NULL, out.df=FALSE)
```

Arguments

| | |
|---------|--|
| formula | A formula of the type $y \sim x_1 + x_2$ where y denotes the name of the categorical variable (a factor in R) which plays the role of the dependent variable, while x_1 and x_2 are the name of the predictors (both categorical variables). Numeric variables are not allowed; eventual numerical variables should be categorized (see function cut) before being passed to <code>pw.assoc</code> . |
| data | The data frame which contains the variables called by <code>formula</code> . |
| weights | The name of the variable in <code>data</code> which provides the units’ weights. Weights are used to estimate frequencies (a cell frequency is estimated by summing the weights of the units which present the given characteristic). Default is <code>NULL</code> (no weights available) and each unit counts 1. When case weight are provided, then they are scales so that their sum equals n , the sample size (assumed to be <code>nrow(data)</code>). |
| out.df | Logical. If <code>NULL</code> measures will be organized in a data frame (a column for each measure). |

Details

This function computes some association, PRE measures, AIC and BIC for each couple response-predictor that can be created starting from argument `formula`. In particular, a two-way contingency table $X \times Y$ is built for each available X variable (X in rows and Y in columns); then the following measures are considered.

Cramer's V :

$$V = \sqrt{\frac{\chi^2}{n \times \min[I - 1, J - 1]}}$$

n is the sample size, I is the number of rows (categories of X) and J is the number of columns (categories of Y). Cramer's V ranges from 0 to 1.

Bias-corrected Cramer's V (V_c) proposed by Bergsma (2013).

Mutual information:

$$I(X; Y) = \sum_{i,j} p_{ij} \log\left(\frac{p_{ij}}{p_{i+}p_{+j}}\right)$$

equal to 0 in case of independence but with infinite upper bound, i.e. $0 \leq I(X; Y) < \infty$. In it $p_{ij} = n_{ij}/n$.

A normalized version of $I(X; Y)$, ranging from 0 (independence) to 1 and not affected by number of categories (I and J):

$$I(X; Y)^* = \frac{I(X; Y)}{\min(H_X, H_Y)}$$

being H_X and H_Y the entropy of the variable X and Y, respectively.

Goodman-Kruskal $\lambda(Y|X)$ (i.e. response conditional on the given predictor):

$$\lambda(Y|X) = \frac{\sum_{i=1}^I \max_j(p_{ij}) - \max_j(p_{+j})}{1 - \max_j(p_{+j})}$$

It ranges from 0 to 1, and denotes how much the knowledge of the row variable X (predictor) helps in reducing the prediction error of the values of the column variable Y (response).

Goodman-Kruskal $\tau(Y|X)$:

$$\tau(Y|X) = \frac{\sum_{i=1}^I \sum_{j=1}^J p_{ij}^2/p_{i+} - \sum_{j=1}^J p_{+j}^2}{1 - \sum_{j=1}^J p_{+j}^2}$$

It takes values in the interval [0,1] and has the same PRE meaning of the lambda.

Theil's uncertainty coefficient:

$$U(Y|X) = \frac{\sum_{i=1}^I \sum_{j=1}^J p_{ij} \log(p_{ij}/p_{i+}) - \sum_{j=1}^J p_{+j} \log p_{+j}}{-\sum_{j=1}^J p_{+j} \log p_{+j}}$$

It takes values in the interval [0,1] and measures the reduction of uncertainty in the column variable Y due to knowing the row variable X. Note that the numerator of $U(Y|X)$ is the mutual information $I(X;Y)$

It is worth noting that λ , τ and U can be viewed as measures of the proportional reduction of the variance of the Y variable when passing from its marginal distribution to its conditional distribution given the predictor X, derived from the general expression (cf. Agresti, 2002, p. 56):

$$\frac{V(Y) - E[V(Y|X)]}{V(Y)}$$

They differ in the way of measuring variance, in fact it does not exist a general accepted definition of the variance for a categorical variable.

Finally, AIC (and BIC) is calculated, as suggested in Sakamoto and Akaike (1977). In particular:

$$AIC(Y|X) = -2 \sum_{i,j} n_{ij} \log \left(\frac{n_{ij}}{n_{i+}} \right) + 2I(J - 1)$$

$$BIC(Y|X) = -2 \sum_{i,j} n_{ij} \log \left(\frac{n_{ij}}{n_{i+}} \right) + I(J - 1) \log(n)$$

being $I(J - 1)$ the parameters (conditional probabilities) to estimate. Note that the R package **catdap** provides functions to identify the best subset of predictors based on AIC.

Please note that the missing values are excluded from the tables and therefore excluded from the estimation of the various measures.

Value

When `out.df=FALSE` (default) a `list` object with four components:

| | |
|----------------------|---|
| <code>V</code> | A vector with the estimated Cramer's V for each couple response-predictor. |
| <code>bcV</code> | A vector with the estimated bias-corrected Cramer's V for each couple response-predictor. |
| <code>mi</code> | A vector with the estimated mutual information $I(X;Y)$ for each couple response-predictor. |
| <code>norm.mi</code> | A vector with the normalized mutual information $I(X;Y)^*$ for each couple response-predictor. |
| <code>lambda</code> | A vector with the values of Goodman-Kruskal $\lambda(Y X)$ for each couple response-predictor. |
| <code>tau</code> | A vector with the values of Goodman-Kruskal $\tau(Y X)$ for each couple response-predictor. |
| <code>U</code> | A vector with the values of Theil's uncertainty coefficient $U(Y X)$ for each couple response-predictor. |
| <code>AIC</code> | A vector with the values of $AIC(Y X)$ for each couple response-predictor. |
| <code>BIC</code> | A vector with the values of $BIC(Y X)$ for each couple response-predictor. |
| <code>npar</code> | A vector with the number of parameters (conditional probabilities) estimated to calculate AIC and BIC for each couple response-predictor. |

When `out.df=TRUE` the output will be a `data.frame` with a column for each measure.

Author(s)

Marcello D’Orazio <mdo.statmatch@gmail.com>

References

- Agresti A (2002) *Categorical Data Analysis. Second Edition*. Wiley, new York.
- Bergsma W (2013) A bias-correction for Cramer’s V and Tschuprow’s T. *Journal of the Korean Statistical Society*, 42, 323–328.
- The Institute of Statistical Mathematics (2018). catdap: Categorical Data Analysis Program Package. R package version 1.3.4. <https://CRAN.R-project.org/package=catdap>
- Sakamoto Y and Akaike, H (1977) Analysis of Cross-Classified Data by AIC. *Ann. Inst. Statist. Math.*, 30, 185-197.

Examples

```
data(quine, package="MASS") #loads quine from MASS
str(quine)

# how Lrn is response variable
pw.assoc(Lrn~Age+Sex+Eth, data=quine)

# usage of units' weights
quine$sww <- runif(nrow(quine), 1,4) #random gen 1<=weights<=4
pw.assoc(Lrn~Age+Sex+Eth, data=quine, weights="sww")
```

RANDwNND.hotdeck

Random Distance hot deck.

Description

This function implements a variant of the distance hot deck method. For each recipient record a subset of of the closest donors is retained and then a donor is selected at random.

Usage

```
RANDwNND.hotdeck(data.rec, data.don, match.vars=NULL,
                  don.class=NULL, dist.fun="Manhattan",
                  cut.don="rot", k=NULL, weight.don=NULL,
                  keep.t=FALSE, ...)
```

Arguments

`data.rec` A numeric matrix or data frame that plays the role of *recipient*. This data frame must contain the variables (columns), specified via `match.vars` and `don.class`, that should be used in the matching.
Missing values (NA) are allowed.

| | |
|-------------------------|---|
| <code>data.don</code> | A matrix or data frame that plays the role of <i>donor</i> . This data frame must contain the variables (columns), specified via <code>match.vars</code> and <code>don.class</code> , that should be used in the matching. |
| <code>match.vars</code> | A character vector with the names of the variables (the columns in both the data frames) that have to be used to compute distances between records (rows) in <code>data.rec</code> and those in <code>data.don</code> . When no matching variables are considered (<code>match.vars=NULL</code>) then all the units in the same donation class are considered as possible donors. Hence one of them is selected at random or with probability proportional to its weight (see argument <code>weight.don</code>). When <code>match.vars=NULL</code> and the donation classes are not created (<code>don.class=NULL</code>) then all the available records in the <code>data.don</code> are considered as potential donors. |
| <code>don.class</code> | A character vector with the names of the variables (columns in both the data frames) that have to be used to identify donation classes. In this case the computation of distances is limited to those units in <code>data.rec</code> and <code>data.doc</code> that belong to the same donation class. The case of empty donation classes should be avoided. It would be preferable that variables used to form donation classes are defined as factor. When not specified (default), no donation classes are used. This may result in a heavy computational effort. |
| <code>dist.fun</code> | A string with the name of the distance function that has to be used. The following distances can be used: “Manhattan” (aka “City block”; default), “Euclidean”, “Mahalanobis”, “exact” or “exact matching”, “Gower”, “minimax”, “difference”, or one of the distance functions available in the package proxy . Note that the distances are computed using the function <code>dist</code> of the package proxy with the exception of the “Gower” (see function <code>gower.dist</code> for details), “Mahalanobis” (function <code>mahalanobis.dist</code>), “minimax” (see <code>maximum.dist</code>) “difference” case. Note that <code>dist.fun="difference"</code> computes just the difference between the values of the unique numeric matching variable considered; in practice, it should be used when the subset of the donation classes should be formed by comparing the values of the unique matching variable (for further details see the argument <code>cut.don</code>). By setting <code>dist.fun="ANN"</code> or <code>dist.fun="RANN"</code> it is possible to search for the <i>k</i> nearest neighbours for each recipient record by using the the Approximate Nearest Neighbor (ANN) search as implemented in the function <code>nn2</code> provided by the package RANN . When <code>dist.fun="Manhattan"</code> , “Euclidean”, “Mahalanobis” or “minimax” all the variables in <code>data.rec</code> and <code>data.don</code> must be numeric. On the contrary, when <code>dist.fun="exact"</code> or <code>dist.fun="exact matching"</code> , all the variables in <code>data.rec</code> and <code>data.don</code> will be converted to character and, as far as the distance computation is concerned, they will be treated as categorical nominal variables, i.e. distance is 0 if a couple of units shows the same response category and 1 otherwise. |
| <code>cut.don</code> | A character string that, jointly with the argument <i>k</i> , identifies the rule to be used to form the subset of the closest donor records. <ul style="list-style-type: none"> • <code>cut.don="rot"</code>: (default) then the number of the closest donors to retain is given by $\lfloor \sqrt{n_D} \rfloor + 1$; being n_D the total number of available donors. In |

this case k must not to be specified.

- `cut.don="span"`: the number of closest donors is determined as the proportion k of all the available donors, i.e. $[n_D \times k]$. Note that, in this case, $0 < k \leq 1$.
- `cut.don="exact"`: the k th closest donors out of the n_D are retained. In this case, $0 < k \leq n_D$.
- `cut.don="min"`: the donors at the minimum distance from the recipient are retained.
- `cut.don="k.dist"`: only the donors whose distance from the recipient is less or equal to the value specified with the argument k . Note that in this case it is not possible to use `dist.fun="ANN"`.
- `cut.don="lt"` or `cut.don("<")`: only the donors whose value of the matching variable is smaller than the value of the recipient are retained. Note that in this case it is has to be set `dist.fun="difference"`.
- `cut.don="le"` or `cut.don("<=")`: only the donors whose value of the matching variable is smaller or equal to the value of the recipient are retained. Note that in this case it is has to be set `dist.fun="difference"`.
- `cut.don="ge"` or `cut.don(">=")`: only the donors whose value of the matching variable is greater or equal to the value of the recipient are retained. Note that in this case it is has to be set `dist.fun="difference"`.
- `cut.don="gt"` or `cut.don(">")`: only the donors whose value of the matching variable is greater than the value of the recipient are retained. Note that in this case it is has to be set `dist.fun="difference"`.

| | |
|-------------------------|--|
| <code>k</code> | Depends on the <code>cut.don</code> argument. |
| <code>weight.don</code> | A character string providing the name of the variable with the weights associated to the donor units in <code>data.don</code> . When this variable is specified, then the selection of a donor among those in the subset of the closest donors is done with probability proportional to its weight (units with larger weight will have a higher chance of being selected). When <code>weight.don=NULL</code> (default) all the units in the subset of the closest donors will have the same probability of being selected. |
| <code>keep.t</code> | Logical, when donation classes are used by setting <code>keep.t=TRUE</code> prints information on the donation classes being processed (by default <code>keep.t=FALSE</code>). |
| <code>...</code> | Additional arguments that may be required by <code>gower.dist</code> , by <code>maximum.dist</code> , by <code>dist</code> or by <code>nn2</code> . |

Details

This function finds a donor record for each record in the recipient data set. The donor is chosen at random in the subset of available donors. This procedure is known as *random hot deck* (cf. Andridge and Little, 2010). In `RANDwNND.hotdeck`, the number of closest donors retained to form the subset is determined according to criterion specified with the argument `cut.don`. The selection of the donor among those in the subset is carried out with equal probability (`weight.don=NULL`) or with probability proportional to a weight associated to the donors, specified via the `weight.don` argument. This procedure is known as *weighted random hot deck* (cf. Andridge and Little, 2010).

The search for the subset of the closest donors can be speed up by using the Approximate Nearest Neighbor search as implemented in the function `nn2` provided by the package `RANN`. Note that this search can be used in all the cases with the exception of `cut.don="k.dist"`.

Note that the same donor can be used more than once.

This function can also be used to impute missing values in a data set. In this case `data.rec` is the part of the initial data set that contains missing values; on the contrary, `data.don` is the part of the data set without missing values. See R code in the Examples for details.

Value

A R list with the following components:

| | |
|-----------------------|---|
| <code>mtc.ids</code> | A matrix with the same number of rows of <code>data.rec</code> and two columns. The first column contains the row names of the <code>data.rec</code> and the second column contains the row names of the corresponding donors selected from the <code>data.don</code> . When the input matrices do not contain row names, then a numeric matrix with the indexes of the rows is provided. |
| <code>sum.dist</code> | A matrix with summary statistics concerning the subset of the closest donors. The first three columns report the minimum, the maximum and the standard deviation of the distances among the recipient record and the donors in the subset of the closest donors, respectively. The 4th column reports the cutting distance, i.e. the value of the distance such that donors at a higher distance are discarded. The 5th column reports the distance between the recipient and the donor chosen at random in the subset of the donors. |
| <code>noad</code> | For each recipient unit, reports the number of donor records in the subset of closest donors. |
| <code>call</code> | How the function has been called. |

Author(s)

Marcello D’Orazio <mdo.statmatch@gmail.com>

References

- Andridge, R.R., and Little, R.J.A. (2010) “A Review of Hot Deck Imputation for Survey Non-response”. *International Statistical Review*, **78**, 40–64.
- D’Orazio, M., Di Zio, M. and Scanu, M. (2006). *Statistical Matching: Theory and Practice*. Wiley, Chichester.
- Rodgers, W.L. (1984). “An evaluation of statistical matching”. *Journal of Business and Economic Statistics*, **2**, 91–102.
- Singh, A.C., Mantel, H., Kinack, M. and Rowe, G. (1993). “Statistical matching: use of auxiliary information as an alternative to the conditional independence assumption”. *Survey Methodology*, **19**, 59–79.

See Also

[NND.hotdeck](#)

Examples

```

data(samp.A, samp.B, package="StatMatch") #loads data sets
?samp.A
?samp.B

# samp.A plays the role of recipient
# samp.B plays the role of donor
# find a donor in the in the same region ("area5") and with the same
# gender ("sex"), then only the closest k=20 donors in terms of
# "age" are considered and one of them is picked up at random

out.RND.1 <- RANDwNND.hotdeck(data.rec=samp.A, data.don=samp.B,
                             don.class=c("area5", "sex"), dist.fun="ANN",
                             match.vars="age", cut.don="exact", k=20)

# create the synthetic (or fused) data.frame:
# fill in "labour5" in A
fused.1 <- create.fused(data.rec=samp.A, data.don=samp.B,
                       mtc.ids=out.RND.1$mtc.ids, z.vars="labour5")
head(fused.1)

# weights ("ww") are used in selecting the donor in the final step

out.RND.2 <- RANDwNND.hotdeck(data.rec=samp.A, data.don=samp.B,
                             don.class=c("area5", "sex"), dist.fun="ANN",
                             match.vars="age", cut.don="exact",
                             k=20, weight.don="ww")
fused.2 <- create.fused(data.rec=samp.A, data.don=samp.B,
                       mtc.ids=out.RND.2$mtc.ids, z.vars="labour5")
head(fused.2)

# find a donor in the in the same region ("area5") and with the same
# gender ("sex"), then only the donors with "age" <= to the age of the
# recipient are considered,
# then one of them is picked up at random

out.RND.3 <- RANDwNND.hotdeck(data.rec=samp.A, data.don=samp.B,
                             don.class=c("area5", "sex"), dist.fun="diff",
                             match.vars="age", cut.don="<=")

# create the synthetic (or fused) data.frame:
# fill in "labour5" in A
fused.3 <- create.fused(data.rec=samp.A, data.don=samp.B,
                       mtc.ids=out.RND.3$mtc.ids, z.vars="labour5")
head(fused.3)

# Example of Imputation of missing values
# introducing missing vales in iris
ir.mat <- iris
miss <- rbinom(nrow(iris), 1, 0.3)
ir.mat[miss==1,"Sepal.Length"] <- NA

```

```

iris.rec <- ir.mat[miss==1,-1]
iris.don <- ir.mat[miss==0,]

#search for NND donors
imp.RND <- RANDwNND.hotdeck(data.rec=iris.rec, data.don=iris.don,
                           match.vars=c("Sepal.Width", "Petal.Length", "Petal.Width"),
                           don.class="Species")

# imputing missing values
iris.rec.imp <- create.fused(data.rec=iris.rec, data.don=iris.don,
                            mtc.ids=imp.RND$mtc.ids, z.vars="Sepal.Length")

# rebuild the imputed data.frame
final <- rbind(iris.rec.imp, iris.don)
head(final)

```

| | |
|-----------------|---------------------------------------|
| rankNND.hotdeck | <i>Rank distance hot deck method.</i> |
|-----------------|---------------------------------------|

Description

This function implements rank hot deck distance method. For each recipient record the closest donors is chosen by considering the distance between the percentage points of the empirical cumulative distribution function.

Usage

```

rankNND.hotdeck(data.rec, data.don, var.rec, var.don=var.rec,
                don.class=NULL, weight.rec=NULL, weight.don=NULL,
                constrained=FALSE, constr.alg="Hungarian",
                keep.t=FALSE)

```

Arguments

| | |
|----------|--|
| data.rec | A numeric matrix or data frame that plays the role of <i>recipient</i> . This data frame must contain the variable <code>var.rec</code> to be used in computing the percentage points of the empirical cumulative distribution function and eventually the variables that identify the donation classes (see argument <code>don.class</code>) and the case weights (see argument <code>weight.rec</code>). Missing values (NA) are not allowed. |
| data.don | A matrix or data frame that plays the role of <i>donor</i> . This data frame must contain the variable <code>var.don</code> to be used in computing percentage points of the the empirical cumulative distribution function and eventually the variables that identify the donation classes (see argument <code>don.class</code>) and the case weights (see argument <code>weight.don</code>). |
| var.rec | A character vector with the name of the variable in <code>data.rec</code> that should be ranked. |

| | |
|--------------------------|--|
| <code>var.rec</code> | A character vector with the name of the variable <code>data.rec</code> that should be ranked. If not specified, by default <code>var.rec=var.rec</code> . |
| <code>don.class</code> | A character vector with the names of the variables (columns in both the data frames) that identify donation classes. In each donation class the computation of percentage points is carried out independently. Then only distances between percentage points of the units in the same donation class are computed. The case of empty donation classes should be avoided. It would be preferable that the variables used to form donation classes are defined as <code>factor</code> . When not specified (default), no donation classes are used. |
| <code>weight.rec</code> | Eventual name of the variable in <code>data.rec</code> that provides the weights that should be used in computing the the empirical cumulative distribution function for <code>var.rec</code> (see Details). |
| <code>weight.don</code> | Eventual name of the variable in <code>data.don</code> that provides the weights that should be used in computing the the empirical cumulative distribution function for <code>var.don</code> (see Details). |
| <code>constrained</code> | Logical. When <code>constrained=FALSE</code> (default) each record in <code>data.don</code> can be used as a donor more than once. On the contrary, when <code>constrained=TRUE</code> each record in <code>data.don</code> can be used as a donor only once. In this case, the set of donors is selected by solving a transportation problem, in order to minimize the overall matching distance. See description of the argument <code>constr.alg</code> for details. |
| <code>constr.alg</code> | A string that has to be specified when <code>constrained=TRUE</code> . Two choices are available: "lpSolve" and "Hungarian". In the first case, <code>constr.alg="lpSolve"</code> , the transportation problem is solved by means of the function <code>lp.transport</code> available in the package <code>lpSolve</code> . When <code>constr.alg="Hungarian"</code> (default) the transportation problem is solved using the Hungarian method, implemented in function <code>solve_LSAP</code> available in the package <code>clue</code> . Note that <code>constr.alg="Hungarian"</code> is faster and more efficient. |
| <code>keep.t</code> | Logical, when donation classes are used by setting <code>keep.t=TRUE</code> prints information on the donation classes being processed (by default <code>keep.t=FALSE</code>). |

Details

This function finds a donor record for each record in the recipient data set. The chosen donor is the one at the closest distance in terms of empirical cumulative distribution (Singh et al., 1990). In practice the distance is computed by considering the estimated empirical cumulative distribution for the reference variable (`var.rec` and `var.don`) in `data.rec` and `data.don`. The empirical cumulative distribution function is estimated by:

$$\hat{F}(y) = \frac{1}{n} \sum_{i=1}^n I(y_i \leq y)$$

being $I() = 1$ if $y_i \leq y$ and 0 otherwise.

In presence of weights, the empirical cumulative distribution function is estimated by:

$$\hat{F}(y) = \frac{\sum_{i=1}^n w_i I(y_i \leq y)}{\sum_{i=1}^n w_i}$$

In the unconstrained case, when there are more donors at the same distance, one of them is chosen at random.

When the donation class are introduced, then the empirical cumulative distribution function is estimated independently in each donation classes and the search of a recipient is restricted to donors in the same donation class.

A donor can be chosen more than once. To avoid it set `constrained=TRUE`. In such a case a donor can be chosen just once and the selection of the donors is carried out by solving a transportation problem with the objective of minimizing the overall matching distance (sum of the distances recipient-donor).

Value

A R list with the following components:

| | |
|----------------------|---|
| <code>mtc.ids</code> | A matrix with the same number of rows of <code>data.rec</code> and two columns. The first column contains the row names of the <code>data.rec</code> and the second column contains the row names of the corresponding donors selected from the <code>data.don</code> . When the input matrices do not contain row names, then a numeric matrix with the indexes of the rows is provided. |
| <code>dist.rd</code> | A vector with the distances between each recipient unit and the corresponding donor. |
| <code>noad</code> | The number of available donors at the minimum distance for each recipient unit (only in unconstrained case) |
| <code>call</code> | How the function has been called. |

Author(s)

Marcello D'Orazio <mdo.statmatch@gmail.com>

References

D'Orazio, M., Di Zio, M. and Scanu, M. (2006). *Statistical Matching: Theory and Practice*. Wiley, Chichester.

Singh, A.C., Mantel, H., Kinack, M. and Rowe, G. (1993). "Statistical matching: use of auxiliary information as an alternative to the conditional independence assumption". *Survey Methodology*, **19**, 59–79.

See Also

[NND.hotdeck](#)

Examples

```
data(samp.A, samp.B, package="StatMatch") #loads data sets

# samp.A plays the role of recipient
?samp.A
```

```

# samp.B plays the role of donor
?samp.B

# rankNND.hotdeck()
# donation classes formed using "area5"
# ecdf computed on "age"
# UNCONSTRAINED case
out.1 <- rankNND.hotdeck(data.rec=samp.A, data.don=samp.B, var.rec="age",
                        don.class="area5")
fused.1 <- create.fused(data.rec=samp.A, data.don=samp.B,
                      mtc.ids=out.1$mtc.ids, z.vars="labour5")
head(fused.1)

# as before but ecdf estimated using weights
# UNCONSTRAINED case
out.2 <- rankNND.hotdeck(data.rec=samp.A, data.don=samp.B, var.rec="age",
                        don.class="area5",
                        weight.rec="ww", weight.don="ww")
fused.2 <- create.fused(data.rec=samp.A, data.don=samp.B,
                      mtc.ids=out.2$mtc.ids, z.vars="labour5")
head(fused.2)

```

| | |
|------------|---|
| rho.bounds | <i>Estimates plausible values of the Pearson's correlation coefficient between two variables observed in distinct samples referred to the same target population.</i> |
|------------|---|

Description

This function assesses the uncertainty in estimating the Pearson's correlation coefficient between $y.rec$ (Y) and $z.don$ (Z) when the two variables are observed in two different samples sharing a number of common predictors.

Usage

```

rho.bounds(data.rec, data.don,
           match.vars, y.rec, z.don,
           w.rec = NULL, w.don = NULL)

```

Arguments

| | |
|----------|--|
| data.rec | dataframe including the Xs (predictors, listed in match.vars) and y.rec (response; target variable in this dataset). |
| data.don | dataframe including the Xs (predictors, listed in match.vars) and z.don (response; target variable in this dataset). |

| | |
|------------|--|
| match.vars | vector with the names of the Xs variables to be used, jointly with y.rec and z.don, in estimating the correlation matrix. If match.vars include one or more factor variables these will be replaced with the corresponding dummies before estimating the correlation matrix. |
| y.rec | character indicating the name of Y target variable in data.rec. It should be a numeric variable. |
| z.don | character indicating the name of Z target variable in data.don. It should be a numeric variable. |
| w.rec | name of the variable with units' weights in data.rec, if available (default NULL); the weights, if provided, are used in estimating the bounds. |
| w.don | name of the variable with units' weights in data.don, if available (default NULL); the weights, if provided, are used in estimating the bounds. |

Details

This function evaluates the uncertainty in the estimation of the Pearson's correlation coefficient between y.rec (Y) and z.don (Z), when the two variables are observed in two different samples that refer to the same target population, but that share a set of common predictors X (match.vars). The evaluation of the uncertainty corresponds to the estimation of the bounds (lower and upper) of the correlation coefficient between Y and Z, given the available data. The method uses the expressions proposed by Rodgers and DeVol (1982). Note that the correlations between the X variables common to both samples (match.vars) are estimated after pooling the samples. Factor variables, if present in match.vars, are replaced by the corresponding dummies before estimating the correlation; this method suffers from a number of critical problems related to the estimation of biserial correlation and the underlying assumption of a Gaussian distribution. The correlation matrix between Y and Xs is estimated on data.rec, while the correlation matrix between Z and Xs is estimated on data.don; this way of working can in some cases give unreliable estimates due to problems with the samples (usually when they are not representative of the same target population).

Value

A vector with three values: the estimated lower bound for Pearson's correlation coefficient between y.rec(Y) and z.don (Z); the estimated upper bound; and, the mid-point of the interval that corresponds to the estimate Pearson's correlation coefficient under the conditional independence assumption (i.e. the correlations between Y and Z is fully explained by the available X variables match.vars).

Author(s)

Marcello D'Orazio <mdo.statmatch@gmail.com>

References

- D'Orazio, M., (2024). *Is Statistical Matching feasible?* Note, https://www.researchgate.net/publication/387699016_Is_statistical_matching_feasible.
- Rodgers, W.L. and DeVol E.B. (1982). An evaluation of statistical matching. *Report Submitted to the Income Survey Development Program*, Dept. of Health and Human Services, Institute for Social Research, University of Michigan.

See Also

[mixed.mtc.](#)

Examples

```
set.seed(11335577)
pos <- sample(x = 1:150, size = 60, replace = FALSE)
ir.A <- iris[pos, c(1:3, 5)]
ir.B <- iris[-pos, c(1:2, 4:5)]

intersect(colnames(ir.A), colnames(ir.B)) # shared Xs

# Xs without Species (factor)
out.1 <- rho.bounds(data.rec=ir.A, data.don=ir.B,
                    match.vars=c("Sepal.Length", "Sepal.Width"),
                    y.rec="Petal.Length", z.don="Petal.Width")
out.1

# Xs with Species (factor)
out.2 <- rho.bounds(data.rec=ir.A, data.don=ir.B,
                    match.vars=c("Sepal.Length", "Sepal.Width", "Species"),
                    y.rec="Petal.Length", z.don="Petal.Width")
out.2
```

| | |
|-----------------|---|
| rho.bounds.pred | <i>Estimates plausible values of the Pearson's correlation coefficient between two variables observed in distinct samples referred to the same target population.</i> |
|-----------------|---|

Description

This function evaluates the uncertainty in estimating the Pearson's correlation coefficient between `y.rec` (Y) and `z.don` (Z) when the two variables are observed in two different samples that share a set of common predictors (Xs). The Xs are used to predict Y and Z respectively, and then the predictions become the input for estimating the uncertainty.

Usage

```
rho.bounds.pred(data.rec, data.don,
                match.vars, y.rec, z.don,
                pred = "lm",
                w.rec = NULL, w.don = NULL,
                out.pred = FALSE, ...)
```

Arguments

| | |
|----------|---|
| data.rec | dataframe including the Xs (predictors, listed in <code>match.vars</code>) and <code>y.rec</code> (response; target variable in this dataset). |
|----------|---|

| | |
|------------|---|
| data.don | dataframe including the Xs (predictors, listed in match.vars) and z.don (response; target variable in this dataset). |
| match.vars | vector with the names of the Xs variables to be used, as (possible) predictors of respectively y.rec and z.don. |
| y.rec | character indicating the name of Y target variable in data.rec. It should be a numeric variable. |
| z.don | character indicating the name of Z target variable in data.don. It should be a numeric variable. |
| pred | String specifying the method used to obtain predictions of both Y and Z. Available methods include pred = "lm" (default) fits two linear regression models (function lm) to get predictions with Y and Z as response variables and match.vars as predictors; pred = "roblm" (default) fits two robust linear regression models (function rlm in package MASS); pred = "lasso" uses the lasso method (R package glmnet, function cv.glmnet) and cross-validation to select a subset of match.vars that are the best predictors of Y (Z) and then obtain the model predictions; pred = "rf" fits randomForest to get predictions of both Y and Z (function randomForest in randomForest). |
| w.rec | possible name of the variable with the weights associated to the units in data.rec, if available; the weights are only used in estimating correlations, not in fitting models. |
| w.don | possible name of the variable with the weights associated to the units in data.don, if available; the weights are only used in estimating correlations, not in fitting models. |
| out.pred | Logical, when TRUE (default is FALSE) the output includes the input datasets with the predictions of both the target variables. |
| ... | addition eventual parameters needed. |

Details

This function evaluates the uncertainty in the estimation of the Pearson's correlation coefficient between y.rec (Y) and z.don (Z), when the two variables are observed in two different samples that refer to the same target population, but that share a set of common predictors X (match.vars). The evaluation of the uncertainty corresponds to the estimation of the bounds (lower and upper) of the correlation coefficient between Y and Z, given the available data. The method uses the expressions proposed by Rodgers and DeVol (1982), but instead of using the Xs match.vars directly, they are replaced by the predictions of both Y and Z provided by the fitted models according to pred. This last way of working avoids the drawbacks encountered when estimating covariances in the presence of several X variables, some of which are categorical (factors) and therefore pose the problem of working with dummies. The final estimation of the bounds is provided by the function rho.bounds. Note that the correlations between the predictions of both Y and Z are estimated after pooling the samples. Survey weights, if available (arguments w.rec and w.don), are used in estimating the correlations, but not in fitting the models.

Value

a list with the following components:

up.rec only when out.pred = TRUE the output list includes data.rec with the predicted values of both Y and Z;

up.don only when out.pred = TRUE the output list includes data.don with the predicted values of both Y and Z;

corr the estimated correlations between Y (Z) and the corresponding predicted values;

bounds a vector with three values: the estimated lower bound for the Pearson's correlation coefficient between y.rec(Y) and z.don (Z); the estimated upper bound; and, the mid-point of the interval that corresponds to the estimate Pearson's correlation coefficient under the conditional independence assumption.

Author(s)

Marcello D'Orazio <mdo.statmatch@gmail.com>

References

D'Orazio, M., (2024). *Is Statistical Matching feasible?* Note, https://www.researchgate.net/publication/387699016_Is_statistical_matching_feasible.

Rodgers, W.L. and DeVol E.B. (1982). An evaluation of statistical matching. *Report Submitted to the Income Survey Development Program*, Dept. of Health and Human Services, Institute for Social Research, University of Michigan.

See Also

[mixed.mtc](#).

Examples

```
set.seed(11335577)
pos <- sample(x = 1:150, size = 60, replace = FALSE)
ir.A <- iris[pos, c(1:3, 5)]
ir.B <- iris[-pos, c(1:2, 4:5)]

intersect(colnames(ir.A), colnames(ir.B)) # shared Xs

op1 <- rho.bounds.pred(data.rec=ir.A, data.don=ir.B,
  match.vars=c("Sepal.Length", "Sepal.Width", "Species"),
  y.rec="Petal.Length", z.don="Petal.Width",
  pred = "lm")
op1
op2 <- rho.bounds.pred(data.rec=ir.A, data.don=ir.B,
  match.vars=c("Sepal.Length", "Sepal.Width", "Species"),
  y.rec="Petal.Length", z.don="Petal.Width",
  pred = "roblm")
op2
```

samp.A

*Artificial data set resembling EU-SILC survey***Description**

This data set provides a limited number of variables observed at persons levels among those usually collected in the European Union Statistics on Income and Living Conditions Survey (EU-SILC). The data are artificially generated, just to show the application of the statistical matching techniques implemented in **StatMatch**.

Usage

```
data(samp.A)
```

Format

A data frame with 3009 observations and the following variables:

HH.P.id unique unit identifier of the type aa.bb where aa identifies the Household while bb identifies the household member

area5 large geographic area, factor with 5 levels: 'NE'=North-East, 'NO'=North-West, 'C'=center, 'S'=South, 'I'=islands

urb Degree of urbanization, factor with 3 levels: '1'=densely populated area, '2'=intermediate area, '3'=thinly populated area

hsize integer, size of the household in which the person lives

hsize5 factor with 5 levels derived from hsize, where the 5th level '>=5' denotes 5 and more people in the household

age integer, the person's age

c.age factor, age categorized in 5 classes

sex factor, the person's gender: '1'=male, '2'=female

marital factor, the person's marital status: '1'=never married, '2'=married, '3'=other (separated, widowed, divorced)

edu7 factor, the person's highest education level attained, follows the ISCED-97 categories: '0'=pre-primary education, '1'=primary education, '2'=lower secondary education, '3'=(upper) secondary education, '4'=post-secondary non tertiary education, '5'=first stage of tertiary education (not leading directly to an advanced research qualification), '6'=second stage of tertiary education (leading to an advanced research qualification)

n.income numeric, the person's net income in Euros

c.neti factor, the person's net income categorized in 7 classes of thousand of Euros

ww numeric, the unit's weight

Details

Please note that this data set is just for illustrative purposes. The unit's weight do not reflect the Italian population size. The variables included are derived starting from the those usually observed in the EU–SILC survey.

Source

This data set is artificially created starting from the EU–SILC survey structure.

References

<https://ec.europa.eu/eurostat/web/income-and-living-conditions/overview>

samp.B

Artificial data set resembling EU–SILC survey

Description

This data set provides a limited number of variables observed at persons levels among those usually collected in the European Union Statistics on Income and Living Conditions Survey (EU–SILC). The data are artificially generated, just to show the application of the statistical matching techniques implemented in **StatMatch**.

Usage

```
data(samp.B)
```

Format

A data frame with 6686 observations and the following variables:

- HH.P.id** unique unit identifier of the type aa.bb where aa identifies the Household while bb identifies the household member
- area5** large geographic area, factor with 5 levels: 'NE'=North–East, 'NO'=North–West, 'C'=center, 'S'=South, 'I'=islands
- urb** Degree of urbanization, factor with 3 levels: '1'=densely populated area, '2'=intermediate area, '3'=thinly populated area
- hsize** integer, size of the household in which the person lives
- hsize5** factor with 5 levels derived from hsize, where the 5th level '>=5' denotes 5 and more people in the household
- age** integer, the person's age
- c.age** factor, age categorized in 5 classes
- sex** factor, the person's gender: '1'=male, '2'=female
- marital** factor, the person's marital status: '1'=never married, '2'=married, '3'=other (separated, widowed, divorced)

edu7 factor, the person's highest education level attained, follows the ISCED-97 categories: '0'=pre-primary education, '1'=primary education, '2'=lower secondary education, '3'=(upper) secondary education, '4'= post-secondary non tertiary education, '5'=first stage of tertiary education (not leading directly to an advanced research qualification), '6'=second stage of tertiary education (leading to an advanced research qualification)

labour5 the person's self-defined economic status, factor with 5 levels: '1'=employee working full-time or part-time, '2'=self-employed working full-time or part-time, '3'=unemployed, '4'=In retirement or in early retirement or has given up business, '5'=other status (student, permanent disabled, in compulsory military service, fulfilling domestic tasks, etc.)

ww numeric, the unit's weight

Details

Please note that this data set is just for illustrative purposes. The unit's weight do not reflect the Italian population size. The variables included are derived starting from the those usually observed in the EU-SILC survey.

Source

This data set is artificially created starting from the EU-SILC survey structure.

References

<https://ec.europa.eu/eurostat/web/income-and-living-conditions/overview>

samp.C

Artificial data set resembling EU-SILC survey

Description

This data set provides a limited number of variables observed at persons levels among those usually collected in the European Union Statistics on Income and Living Conditions Survey (EU-SILC). The data are artificially generated, just to show the application of the statistical matching techniques implemented in **StatMatch**.

Usage

`data(samp.C)`

Format

A data frame with 980 observations and the following variables:

HH.P.id unique unit identifier of the type aa.bb where aa identifies the Household while bb identifies the household member

area5 large geographic area, factor with 5 levels: 'NE'=North-East, 'NO'=North-West, 'C'=center, 'S'=South, 'I'=islands

urb Degree of urbanization, factor with 3 levels: ‘1’=densely populated area, ‘2’=intermediate area, ‘3’=thinly populated area

hsize integer, size of the household in which the person lives

hsize5 factor with 5 levels derived from hsize, where the 5th level ‘>=5’ denotes 5 and more people in the household

age integer, the person’s age

c.age factor, age categorized in 5 classes

sex factor, the person’s gender: ‘1’=male, ‘2’=female

marital factor, the person’s marital status: ‘1’=never married, ‘2’=married, ‘3’=other (separated, widowed, divorced)

edu7 factor, the person’s highest education level attained, follows the ISCED-97 categories: ‘0’=pre-primary education, ‘1’=primary education, ‘2’=lower secondary education, ‘3’= (upper) secondary education, ‘4’= post-secondary non tertiary education, ‘5’=first stage of tertiary education (not leading directly to an advanced research qualification), ‘6’=second stage of tertiary education (leading to an advanced research qualification)

labour5 the person’s self-defined economic status, factor with 5 levels: ‘1’=employee working full-time or part-time, ‘2’=self-employed working full-time or part-time, ‘3’=unemployed, ‘4’=In retirement or in early retirement or has given up business, ‘5’=other status (student, permanent disabled, in compulsory military service, fulfilling domestic tasks, etc.)

n.income numeric, the person’s net income in Euros

c.neti factor, the person’s net income categorized in 7 classes of thousand of Euros

ww numeric, the unit’s weight

Details

Please note that this data set is just for illustrative purposes. The unit’s weight do not reflect the Italian population size. The variables included are derived starting from the those usually observed in the EU-SILC survey.

Source

This data set is artificially created starting from the EU-SILC survey structure.

References

<https://ec.europa.eu/eurostat/web/income-and-living-conditions/overview>

| | |
|---------------|---|
| selMtc.by.unc | <i>Identifies the best combination if matching variables in reducing uncertainty in estimation the contingency table Y vs. Z.</i> |
|---------------|---|

Description

This function identifies the “best” subset of matching variables in terms of reduction of uncertainty when estimating relative frequencies in the contingency table Y vs. Z. The sequential procedure presented in D’Orazio *et al.* (2017 and 2019) is implemented. This procedure avoids exploring all the possible combinations of the available X variables as in [Fbwidths.by.x](#).

Usage

```
selMtc.by.unc(tab.x, tab.xy, tab.xz, corr.d=2,
              nA=NULL, nB=NULL, align.margins=FALSE)
```

Arguments

| | |
|--------|---|
| tab.x | A R table crossing the X variables. This table must be obtained by using the function xtabs or table , e.g. <code>tab.x <- xtabs(~x1+x2+x3, data=data.all)</code> . A minimum number of 3 variables is needed. |
| tab.xy | A R table of X vs. Y variable. This table must be obtained by using the function xtabs or table , e.g. <code>table.xy <- xtabs(~x1+x2+x3+y, data=data.A)</code> . A single categorical Y variables is allowed. At least three categorical variables should be considered as X variables (common variables). The same X variables in tab.x must be available in tab.xy. Usually, it is assumed that the joint distribution of the X variables computed from tab.xy is equal to tab.x (a warning appears if any absolute difference is greater than tol). Note that when the marginal distribution of X in tab.xy is not equal to that of tab.x it is possible to align them before computations (see argument align.margins). |
| tab.xz | A R table of X vs. Z variable. This table must be obtained by using the function xtabs or table , e.g. <code>tab.xz <- xtabs(~x1+x2+x3+z, data=data.B)</code> . A single categorical Z variable is allowed. At least three categorical variables should be considered as X variables (common variables). The same X variables in tab.x must be available in tab.xz. Usually, it is assumed that the joint distribution of the X variables computed from tab.xz is equal to tab.x (a warning appears if any absolute difference is greater than tol). Note that when the marginal distribution of X in tab.xz is not equal to that of tab.x it is possible to align them before computations (see argument align.margins). |
| corr.d | Integer, indicates the penalty that should be introduced in estimating the uncertainty by means of the average width of cell bounds. When corr.d=1 the penalty being considered is the one introduced in D’Orazio <i>et al.</i> (2017) (i.e. penalty1 in Fbwidths.by.x). When corr.d=2 (default) it is considered a penalty suggested |

in D’Orazio *et al.* (2019) (indicated as “penalty2” in `Fbwidths.by.x`). Finally, no penalties are considered when `corr.d=0`.

| | |
|---------------|--|
| nA | Integer, sample size of file A used to estimate <code>tab.xy</code> . If NULL is obtained as sum of frequencies in <code>tab.xy</code> . |
| nB | Integer, sample size of file B used to estimate <code>tab.xz</code> . If NULL is obtained as sum of frequencies in <code>tab.xz</code> . |
| align.margins | Logical (default FALSE). When when TRUE the distribution of X variables in <code>tab.xy</code> is aligned with the distribution resulting from <code>tab.x</code> , without affecting the marginal distribution of Y. Similarly the distribution of X variables in <code>tab.xz</code> is aligned with the distribution resulting from <code>tab.x</code> , without affecting the marginal distribution of Z. The alignment is performed by running IPF algorithm as implemented in the function <code>Estimate</code> in the package <code>mipfp</code> . To avoid lack of convergence due to combinations of Xs encountered in one table but not in the other (statistical 0s), before running IPF a small constant (1e-06) is added to empty cells in <code>tab.xy</code> and <code>tab.xz</code> . |

Details

This function follows the sequential procedure described in D’Orazio *et al.* (2017, 2019) to identify the combination of common variables most effective in reducing uncertainty when estimating the contingency table Y vs. Z. Initially, the available Xs are ordered according to the reduction of average width of uncertainty bounds when conditioning on each of them. Then in each step one the remaining X variables is added until the table became too sparse; in practice the procedure stops when:

$$\min \left[\frac{n_A}{H_{D_m} \times J}, \frac{n_B}{H_{D_m} \times K} \right] \leq 1$$

For major details see also `Fbwidths.by.x`.

Value

A list with the main outcomes of the procedure.

| | |
|---------|--|
| ini.ord | Average width of uncertainty bounds when conditioning on each of the available X variables. Variable most effective in reducing uncertainty comes first. The ordering determines the order in which they are entered in the sequential procedure. |
| list.xs | List with the various combinations of the matching variables being considered in each step. |
| av.df | Data.frame with all the relevant information for each of combination of X variables. The last row corresponds to the combination of the X variables identified as the best in reducing average width of uncertainty bounds (penalized or not depending on the input argument <code>corr.d</code>). For each combination of X variables the following additional information are reported: the number of cells (name starts with “nc”); the number of empty cells (name starts with “nc0”); the average relative frequency (name starts with “av.crf”); sparseness measured as Cohen’s |

effect size with respect to equiprobability (uniform distribution across cells). Finally there are the value of the stopping criterion (“min.av”), the unconditioned average width of uncertainty bounds (“avw”), the penalty term (“penalty”) and the penalized width (“avw.pen”; $avw.pen=avw+penalty$).

Author(s)

Marcello D’Orazio <mdo.statmatch@gmail.com>

References

- D’Orazio, M., Di Zio, M. and Scanu, M. (2006). *Statistical Matching: Theory and Practice*. Wiley, Chichester.
- D’Orazio, M., Di Zio, M. and Scanu, M. (2017). “The use of uncertainty to choose matching variables in statistical matching”. *International Journal of Approximate Reasoning*, 90, pp. 433-440.
- D’Orazio, M., Di Zio, M. and Scanu, M. (2019). “Auxiliary variable selection in a statistical matching problem”. In Zhang, L.-C. and Chambers, R. L. (eds.) *Analysis of Integrated Data*, Chapman & Hall/CRC (forthcoming).

See Also

[Fbwidths.by.x](#), [Frechet.bounds.cat](#)

Examples

```
data(quine, package="MASS") #loads quine from MASS
str(quine)
quine$c.Days <- cut(quine$Days, c(-1, seq(0,50,10),100))
table(quine$c.Days)

# split quine in two subsets
suppressWarnings(RNGversion("3.5.0"))
set.seed(1111)
lab.A <- sample(nrow(quine), 70, replace=TRUE)
quine.A <- quine[lab.A, 1:4]
quine.B <- quine[-lab.A, c(1:3,6)]

# compute the tables required by Fbwidths.by.x()
freq.xA <- xtabs(~Eth+Sex+Age, data=quine.A)
freq.xB <- xtabs(~Eth+Sex+Age, data=quine.B)

freq.xy <- xtabs(~Eth+Sex+Age+Lrn, data=quine.A)
freq.xz <- xtabs(~Eth+Sex+Age+c.Days, data=quine.B)

# apply Fbwidths.by.x()
bb <- Fbwidths.by.x(tab.x=freq.xA+freq.xB,
                   tab.xy=freq.xy, tab.xz=freq.xz,
                   warn=FALSE)

bb$sum.unc
```

```
cc <- selMtc.by.unc(tab.x=freq.xA+freq.xB,  
                   tab.xy=freq.xy, tab.xz=freq.xz, corr.d=0)  
cc$ini.ord  
cc$av.df
```

Index

- * **cluster**
 - fact2dummy, 17
 - gower.dist, 28
- * **datasets**
 - samp.A, 71
 - samp.B, 72
 - samp.C, 73
- * **dplot**
 - plotBounds, 51
 - plotCont, 52
 - plotTab, 54
- * **manip**
 - create.fused, 13
 - create.imputed, 15
- * **multivariate**
 - comp.cont, 8
 - comp.prop, 10
 - fact2dummy, 17
 - Fbwidths.by.x, 21
 - Frechet.bounds.cat, 24
 - gower.dist, 28
 - mahalanobis.dist, 36
 - maximum.dist, 37
 - pBayes, 48
 - plotBounds, 51
 - plotCont, 52
 - plotTab, 54
 - pw.assoc, 55
 - selMtc.by.unc, 75
- * **nonparametric**
 - mixed.mtc, 39
 - NND.hotdeck, 44
 - RANDwNND.hotdeck, 58
 - rankNND.hotdeck, 63
- * **regression**
 - mixed.mtc, 39
- * **survey**
 - comb.samples, 3
 - harmonize.x, 31
 - calibrate, 4–6, 32, 33
 - comb.samples, 3, 33
 - comp.cont, 8, 54
 - comp.prop, 10, 55
 - cov.rob, 36
 - create.fused, 13
 - create.imputed, 15
 - cut, 55
 - cv.glmnet, 19, 69
 - daisy, 30
 - dist, 30, 45, 59, 60
 - Estimate, 25, 76
 - fact2dummy, 17
 - Fbounds.pred, 18
 - Fbwidths.by.x, 21, 27, 75–77
 - formula, 4, 32
 - Frechet.bounds.cat, 19, 20, 22, 24, 24, 51, 77
 - gower.dist, 18, 28, 45, 59, 60
 - harmonize.x, 5, 6, 24, 27, 31
 - lp.transport, 40, 45, 64
 - mahalanobis, 37
 - mahalanobis.dist, 36, 42, 45, 59
 - maximum.dist, 37, 45, 59, 60
 - mixed.mtc, 14, 16, 39, 68, 70
 - model.matrix, 4, 32
 - multinom, 19
 - naive_bayes, 19
 - nclass, 9, 53
 - nn2, 59, 60
 - NND.hotdeck, 14, 16, 42, 44, 61, 65
 - pBayes, 26, 48

plotBounds, [51](#)
plotCont, [10](#), [52](#)
plotTab, [54](#)
postStratify, [32](#), [33](#)
pw.assoc, [55](#)

randomForest, [19](#), [69](#)
RANDwNND.hotdeck, [14](#), [16](#), [47](#), [58](#)
rank, [30](#), [38](#)
rankNND.hotdeck, [14](#), [16](#), [63](#)
rho.bounds, [66](#), [69](#)
rho.bounds.pred, [68](#)

samp.A, [71](#)
samp.B, [72](#)
samp.C, [73](#)
selMtc.by.unc, [75](#)
solve_LSAP, [40](#), [45](#), [64](#)
StatMatch (StatMatch-package), [2](#)
StatMatch-package, [2](#)
svydesign, [3](#), [6](#), [31](#), [33](#)

table, [10](#), [21](#), [22](#), [25](#), [48](#), [75](#)

wtd.quantile, [9](#), [53](#)

xtabs, [10](#), [21](#), [22](#), [25](#), [48](#), [75](#)