

Package ‘rcbalance’

October 14, 2022

Type Package

Title Large, Sparse Optimal Matching with Refined Covariate Balance

Version 1.8.8

Date 2022-3-25

Author Samuel D. Pimentel

Maintainer Samuel D. Pimentel <spi@berkeley.edu>

Description Tools for large, sparse optimal matching of treated units and control units in observational studies. Provisions are made for refined covariate balance constraints, which include fine and near-fine balance as special cases. Matches are optimal in the sense that they are computed as solutions to network optimization problems rather than greedy algorithms. See Pimentel, et al.(2015) <doi:10.1080/01621459.2014.997879> and Pimentel (2016), Obs. Studies 2(1):4-23. The rrelaxiv package, which provides an alternative solver for the underlying network flow problems, carries an academic license and is not available on CRAN, but may be downloaded from Github at <<https://github.com/josherrickson/rrelaxiv/>>.

Depends R (>= 3.2.0), MASS, plyr

Imports rlemon

License MIT + file LICENSE

Suggests optmatch, testthat, rrelaxiv

Additional_repositories <https://errickson.net/rrelaxiv/>

NeedsCompilation no

Repository CRAN

Date/Publication 2022-03-25 23:30:02 UTC

R topics documented:

rcbalance-package 2

build.dist.struct	3
callrelax	5
count.pairings	6
dist2net	6
rcbalance	8

Index	12
--------------	-----------

rcbalance-package	<i>Large, Sparse Optimal Matching with Refined Covariate Balance</i>
-------------------	--

Description

Tools for large, sparse optimal matching of treated units and control units in observational studies. Provisions are made for refined covariate balance constraints, which include fine and near-fine balance as special cases. Matches are optimal in the sense that they are computed as solutions to network optimization problems rather than greedy algorithms. See Pimentel, et al.(2015) <doi:10.1080/01621459.2014.997879> and Pimentel (2016), *Obs. Studies* 2(1):4-23. The rrelaxiv package, which provides an alternative solver for the underlying network flow problems, carries an academic license and is not available on CRAN, but may be downloaded from Github at <<https://github.com/josherrickson/rrelaxiv/>>.

Details

The DESCRIPTION file:

```
Package:          rcbalance
Type:            Package
Title:           Large, Sparse Optimal Matching with Refined Covariate Balance
Version:         1.8.8
Date:           2022-3-25
Author:          Samuel D. Pimentel
Maintainer:     Samuel D. Pimentel <spi@berkeley.edu>
Description:    Tools for large, sparse optimal matching of treated units and control units in observational studies.
Depends:        R (>= 3.2.0), MASS, plyr
Imports:        rlemon
License:        MIT + file LICENSE
Suggests:       optmatch, testthat, rrelaxiv
Additional_repositories: https://errickson.net/rrelaxiv/
```

Index of help topics:

```
build.dist.struct  Build Distance Structure for Matching with
                   Refined Balance
callrelax          Solve Network Flow Problem using External
                   Solver
count.pairings     Count treatment-control pairings.
```

dist2net	Building and Manipulating Network Flow Problems
rcbalance	Optimal Matching with Refined Covariate Balance
rcbalance-package	Large, Sparse Optimal Matching with Refined Covariate Balance

This package computes sparse matches that are optimal under a set of refined covariate balance constraints. These constraints, provided by the user, are a set of nested categorical variables of decreasing importance which must be marginally balanced as closely as possible in the resulting treated and matched control populations. For more detail see the references.

The main function is `rcbalance`, which takes a distance/sparsity object containing information about matchability of the treated and control units and a list of fine balance variables and produces a match. The `build.dist.struct` function can be used to construct the distance/sparsity object from covariate information. The `count.pairings` function can be used to assess the sparsity of a proposed match. The other functions are largely for internal use and should not be needed by the large majority of users.

By default the package uses the R package `rlemon` to solve the minimum-cost network flow optimization problems by which matches are computed. Alternatively, users may specify that the `rrelaxiv` package should be used instead. However, this package carries an academic license and is not available on CRAN so users must install it themselves.

Author(s)

Samuel D. Pimentel

Maintainer: Samuel D. Pimentel <spi@berkeley.edu>

References

Pimentel, S.D., Kelz, R.R., Silber, J.H., and Rosenbaum, P.R. (2015) Large, sparse optimal matching with refined covariate balance in an observational study of the health outcomes produced by new surgeons, *JASA* 110 (510), 515-527.

Pimentel, S.D. (2016) Large, sparse optimal matching with R package `rcbalance`, *Obs. Studies* 2, 4-23.

build.dist.struct	<i>Build Distance Structure for Matching with Refined Balance</i>
-------------------	---

Description

This function computes rank-based Mahalanobis distances between treated and control units and returns an object suitable for use in the `distance.structure` argument of `rcbalance`.

Usage

```
build.dist.struct(z, X, exact = NULL, calip.option = "propensity",
calip.cov = NULL, caliper = 0.2, verbose = FALSE)
```

Arguments

z	a vector of treatment and control indicators, 1 for treatment and 0 for control.
X	a data frame or a numeric or logical matrix containing covariate information for treated and control units. Its row count must be equal to the length of z.
exact	an optional vector of the same length as z. If this argument is specified, treated units will only be allowed to match to control units that have equal values in the corresponding indices of the exact vector. For example, to match patients within hospitals only, one could set exact equal to a vector of hospital IDs for each patient.
calip.option	one of ('propensity', 'user', 'none'). If 'propensity' is specified (the default option), the function estimates a propensity score via logistic regression of z on X and imposes a propensity score caliper. If 'user' is specified, the user must provide a vector of values on which a caliper will be enforced using the calip.cov argument. If 'none' is specified no caliper is used.
calip.cov	see calip.option.
caliper	gives the size of the caliper when the user specifies the calip.option argument as 'propensity' or 'calip.cov'.
verbose	if TRUE, prints output describing specific adjustments made in creating the distance objects.

Details

If X is a data frame and contains any character variables they are converted to factors with a warning. If there are missing values in factor columns of X, they are treated as a new factor level. If there are missing values in numeric or logical columns of X, an indicator of missingness for that column is added to X and the missing values are imputed with the column mean. This follows the recommendations of Rosenbaum (*Design of Observational Studies* section 9.4, 2010).

Value

A distance.structure object, the form of which is described in the documentation for the distance.structure argument of rcbalance. Treated and control indices are numbered 1:nt and 1:nc respectively based on the order in which they appear in the z vector.

Author(s)

Samuel D. Pimentel

See Also

[rcbalance](#)

 callrelax

Solve Network Flow Problem using External Solver

Description

Solves network flow optimization problems by calling an external solver, either the Lemon Optimization library or the RELAX-IV algorithm, as implemented in FORTRAN by Dimitri Bertsekas and Paul Tseng.

IMPORTANT NOTE 1: the RELAX-IV code is bound by an academic license; as a result the rrelaxiv package implementing it is not available on CRAN and must be downloaded separately.

Usage

```
callrelax(net, solver = 'rlemon')
```

Arguments

net	<p>a network flow problem, formatted as a list with the following arguments (where the network contains nnode nodes, numbered 1 through nnode and narc arcs):</p> <ul style="list-style-type: none"> • startn: a vector of length narc containing the node numbers of the start nodes of each arc in the network. • endn: a vector of length narc containing the node numbers of the end nodes of each arc in the network. • ucap: a vector of length narc containing the (integer) upper capacity of each arc in the network. • cost: a vector of length narc containing the (integer) cost of each arc in the network. • b: a vector of length nnode containing the (integer) supply or demand of each node in the network. Supplies are given as positive numbers and demands as negative numbers.
solver	<p>the name of the package used to solve the network flow optimization problem underlying the match, one of 'rlemon' (which uses the Lemon Optimization Library) and 'rrelaxiv' (which uses the RELAX-IV algorithm).</p>

Value

A list with the following elements:

crash	an integer equal to zero (included in order to support legacy versions of code).
feasible	an integer, equal to zero if the problem is infeasible and equal to 1 if it is feasible. A network with a supply/demand vector of all zeroes may also be flagged as infeasible.
x	a vector equal in length to the number of arcs in argument problem net, giving in each coordinate the number of units of flow passing across the corresponding edge in the optimal network flow.

Author(s)

Samuel D. Pimentel

count.pairings	<i>Count treatment-control pairings.</i>
----------------	--

Description

Given a treatment indicator and a potential blocking variable, counts the number of allowed treatment-control pairings in the whole match within blocks of the proposed variable.

Usage

```
count.pairings(z, exact)
```

Arguments

z	a vector of treatment indicators. Must contain exactly 2 distinct values, one for treated and one for control.
exact	a vector of categories of a potential blocking variable. Must be the same length as argument z.

Value

The number of within-block treatment-control edges contained in the sparse match with the proposed blocks.

Author(s)

Samuel D. Pimentel

dist2net	<i>Building and Manipulating Network Flow Problems</i>
----------	--

Description

These are internal rcbalance methods not meant to be called directly by users. They are used to construct a network flow problem from the information about a matching problem that is passed to the rcbalance method.

Usage

```
dist2net(dist.struct, k, exclude.treated = FALSE, ncontrol = NULL)
```

```
dist2net.matrix(dist.struct, k, exclude.treated = FALSE)
```

```
add.layer(net.layers, new.layer)
```

```
penalty.update(net.layers, newtheta, newp = NA)
```

```
penalize.near.exact(net.layers, near.exact)
```

Arguments

<code>dist.struct</code>	An object specifying the sparsity structure of the match. For the <code>dist2net</code> method it is a list of vectors, and for the <code>dist2net.matrix</code> method it is a matrix or <code>InfinitySparseMatrix</code> . See <code>rcbalance</code> documentation for more details.
<code>k</code>	a nonnegative integer. The number of control units to which each treated unit will be matched.
<code>exclude.treated</code>	if TRUE, then when there is no feasible match using all treated units, a minimal number of treated units may be dropped so that a match can be formed. Specifying this argument adds penalized edges to the network so that such a match can be computed. NOTE: this argument is incompatible with values of <code>k</code> greater than 1.
<code>ncontrol</code>	the number of controls in the matching problem. If left NULL (the default value), the value will be intuited from the maximum control label in the sparsity object.
<code>net.layers</code>	a layered network object of the type produced by the <code>dist2net</code> function.
<code>new.layer</code>	a vector equal in length to the number of treated and control units in the matching problem. Each coordinate contains the value of a new fine balance variable for the corresponding unit.
<code>newtheta</code>	optional argument giving a new value for the <code>theta</code> field of the <code>net.layers</code> object (see value section for description of this field).
<code>newp</code>	optional argument giving a new value for the <code>p</code> field of the <code>net.layers</code> object (see value section for description of this field).
<code>near.exact</code>	a vector equal in length to the number of treated and control units in the matching problem. Edges between units with different values of this variable will be penalized.

Details

`dist2net` and `dist2net.matrix` take the distance structure given to `rcbalance` encoding information about the matching problem and converts it into a network flow problem. `add.layer` adds network structure to handle an individual fine balance variable (it can be called iteratively to add many such variables). `penalty.update` is used to change the penalties for each layer (and the penalties for edges used to exclude treated units if they are present) and `penalize.near.exact` is used to add penalties to the treated-control edges to allow near-exact matching. See the references for a detailed description of how the matching problem is transformed into a network.

Value

A layered network object, formatted as a list with the following arguments (where `narcs` is the number of arcs and `nnode` is the number of nodes in the network):

<code>startn</code>	a vector of length <code>narc</code> containing the node numbers of the start nodes of each arc in the network.
<code>endn</code>	a vector of length <code>narc</code> containing the node numbers of the end nodes of each arc in the network.
<code>ucap</code>	a vector of length <code>narc</code> containing the (integer) upper capacity of each arc in the network.
<code>cost</code>	a vector of length <code>narc</code> containing the (integer) cost of each arc in the network.
<code>b</code>	a vector of length <code>nnode</code> containing the (integer) supply or demand of each node in the network. Supplies are given as positive numbers and demands as negative numbers.
<code>tcarcs</code>	an integer giving the total number of arcs between the treated and control nodes in the network.
<code>layers</code>	a list object containing information about the refined covariate balance layers of the network.
<code>z</code>	a vector of treatment indicators.
<code>fb.structure</code>	a matrix containing information about the membership of the treated and control units in the different classes of refined balance covariates.
<code>penalties</code>	a vector of integer penalties, one for each fine balance layer.
<code>theta</code>	a value no less than 1 giving the ratio by which the penalty is increased with each additional layer of fine balance.
<code>p</code>	a nonnegative value giving the penalty for the finest level of fine balance.

Author(s)

Samuel D. Pimentel

rcbalance

Optimal Matching with Refined Covariate Balance

Description

This function computes an optimal match with refined covariate balance.

Usage

```
rcbalance(distance.structure, near.exact = NULL, fb.list = NULL,
treated.info = NULL, control.info = NULL, exclude.treated = FALSE, target.group = NULL,
k = 1, penalty = 3, tol = 1e-5, solver = 'rlemon')
```


Arguments

<code>distance.structure</code>	<p>a list of vectors that encodes information about covariate distances between treated and control units. The list is equal in length to the number of treated units. Each vector corresponds to a treated unit and is equal in length to the number of control units to which it can be matched. It is assumed that there are a total of <code>nc</code> control units in the problem and that they are numbered from 1 to <code>nc</code>. The names of each vector in the list give the index (in the vector <code>1:nc</code>) of the control units to which the treated unit in question can be matched, and the elements of each vector are the covariate distances between the treated unit and the corresponding control. Note that for a dense matching problem (in which each treated unit can be matched to any control), every vector in the list will have length <code>nc</code> and rownames 1 through <code>nc</code>.</p> <p>Alternatively, this same information can be passed as a <code>matrix</code> or <code>InfinitySparseMatrix</code> with rows corresponding to treated units and columns corresponding to controls. Entries given as <code>Inf</code> correspond to pairs that cannot be matched.</p>
<code>near.exact</code>	<p>an optional character vector specifying names of covariates for near-exact matching. This argument takes precedence over any refined covariate balance constraints, so the match will produce the best refined covariate balance subject to matching exactly on this variable wherever possible. If multiple covariates are named, near-exact matching will be done on their interaction.</p>
<code>fb.list</code>	<p>an optional list of character vectors specifying covariates to be used for refined balance. Each element of the list corresponds to a level of refined covariate balance, and the levels are assumed to be in decreasing order of priority. Each character vector should contain one or more names of categorical covariates on which the user would like to enforce near fine balance. If multiple covariates are specified, an interaction is created between the categories of the covariates and near fine balance is enforced on the interaction. IMPORTANT: covariates or interactions coming later in the list must be nested within covariates coming earlier in the list; if this is not the case the function will stop with an error. An easy way to ensure that this occurs is to include in each character vector all the variables named in earlier list elements. If the <code>fb.list</code> argument is specified, the <code>treated.info</code> and <code>control.info</code> arguments must also be specified.</p>
<code>treated.info</code>	<p>an optional data frame containing covariate information for the treated units in the problem. The row count of this data frame must be equal to the length of the <code>distance.structure</code> argument, and it is assumed that row <code>i</code> contains covariate information for the treated unit described by element <code>i</code> of <code>distance.structure</code>. In addition, the column count and column names must be identical to those of the <code>control.info</code> argument, and the column names must include all of the covariate names mentioned in the <code>near.exact</code> and <code>fb.list</code> arguments.</p>
<code>control.info</code>	<p>an optional data frame containing covariate information for the control units in the problem. The row count of this data frame must be no smaller than the maximum control index in the <code>distance.structure</code> argument, and it is assumed that row <code>i</code> contains the covariate information for the control indexed by <code>i</code> in <code>distance.structure</code>. In addition, the column count and column names must be identical to those of the <code>treated.info</code> argument.</p>

<code>exclude.treated</code>	if TRUE, then when there is no feasible match using all treated units, a minimal number of treated units will be dropped so that a match can be formed. The excluded treated units will be selected optimally so that the cost of the matching is reduced as much as possible. NOTE: <code>exclude.treated = TRUE</code> is incompatible with arguments to <code>target.group</code> and with values of <code>k</code> larger than 1.
<code>target.group</code>	an optional data frame of observations with the desired covariate distribution for the selected control group, if it differs from the covariate distribution of the treated units. This argument will be ignored unless <code>fb.list</code> , <code>treated.info</code> and <code>control.info</code> are also specified, and it must have the same dimensions as <code>treated.info</code> .
<code>k</code>	a nonnegative integer. The number of control units to which each treated unit will be matched.
<code>penalty</code>	a value greater than 1. This is a tuning parameter that helps ensure the different levels of refined covariate balance are prioritized correctly. Setting the penalty higher tends to improve the guarantee of match optimality up to a point, but penalties above a certain level cause integer overflows and throw errors. Usually it is not recommended that the user change this parameter from its default value.
<code>tol</code>	edge cost tolerance. This is the smallest tolerated difference between matching costs; cost differences smaller than this will be considered zero. Match distances will be scaled by inverse tolerance, so when matching with large edge costs or penalties the tolerance may need to be increased.
<code>solver</code>	the name of the package used to solve the network flow optimization problem underlying the match, one of 'rlemon' (which uses the Lemon Optimization Library) and 'rrelaxiv' (which uses the RELAX-IV algorithm).

Details

To use the option `solver = 'rrelaxiv'`, the user must install the `rrelaxiv` manually; it is not hosted on CRAN because it carries an academic license.

Value

A list with the following components:

<code>matches</code>	a <code>nt</code> by <code>k</code> matrix containing the matched sets produced by the algorithm (where <code>nt</code> is the number of treated units). The rownames of this matrix are the numbers of the treated units (indexed by their position in <code>distance.structure</code>), and the elements of each row contain the indices of the control units to which this treated unit has been matched.
<code>fb.tables</code>	a list of matrices, equal in length to the <code>fb.list</code> argument. Each matrix is a contingency table giving the counts among treated units and matched controls for each level of the categorical variable specified by the corresponding element of <code>fb.list</code> .

Author(s)

Samuel D. Pimentel

References

Pimentel, S.D., Kelz, R.R., Silber, J.H., and Rosenbaum, P.R. (2015) Large, sparse optimal matching with refined covariate balance in an observational study of the health outcomes produced by new surgeons, *JASA* 110 (510), 515-527.

Examples

```
## Not run:
library(optmatch)
data(nuclearplants)

#require exact match on variables ne and pt, use rank-based Mahalanobis distance
my.dist.struct <- build.dist.struct(z = nuclearplants$pr,
X = subset(nuclearplants[c('date', 't1', 't2', 'cap', 'bw', 'cum.n')]),
exact = paste(nuclearplants$ne, nuclearplants$pt, sep = '.'))

#match with refined covariate balance, first on ct then on (ct x bw)
rcbalance(my.dist.struct, fb.list = list('ct',c('ct','bw')),
  treated.info = nuclearplants[which(nuclearplants$pr ==1),],
  control.info = nuclearplants[which(nuclearplants$pr == 0),])

#repeat the same match using match_on tool from optmatch and regular Mahalanobis distance
exact.mask <- exactMatch(pr ~ ne + pt, data = nuclearplants)
my.dist.matrix <- match_on(pr ~ date + t1 + t2 + cap + bw + cum.n,
within = exact.mask, data = nuclearplants)
match.matrix <-
rcbalance(my.dist.matrix*100, fb.list = list('ct',c('ct','bw')),
  treated.info = nuclearplants[which(nuclearplants$pr ==1),],
  control.info = nuclearplants[which(nuclearplants$pr == 0),])

## End(Not run)
```

Index

`add.layer (dist2net)`, 6

`build.dist.struct`, 3

`callrelax`, 5

`count.pairings`, 6

`dist2net`, 6

`penalize.near.exact (dist2net)`, 6

`penalty.update (dist2net)`, 6

`rcbalance`, 4, 8

`rcbalance-package`, 2

`remove.layer (dist2net)`, 6