

Package ‘rcsubset’

November 8, 2021

Type Package

Title Optimal Subset Matching with Refined Covariate Balance

Version 1.1.5

Date 2021-10-15

Author Samuel D. Pimentel

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

Description Tools for optimal subset matching of treated units and control units in observational studies, with support for refined covariate balance constraints, (including fine and near-fine balance as special cases). A close relative is the 'rcbalance' package. See Pimentel, et al.(2015) <doi:10.1080/01621459.2014.997879> and Pimentel and Kelz (2020) <doi:10.1080/01621459.2020.1720693>.

Depends R (>= 3.2.0), MASS, plyr

Imports rcbalance

License MIT + file LICENSE

Suggests optmatch, testthat

NeedsCompilation no

Repository CRAN

Date/Publication 2021-11-08 22:50:02 UTC

R topics documented:

rcsubset-package	2
callrelax	3
dist2net	4
rcsubset	6

Index	10
--------------	-----------

rcbsubset-package

*Optimal Subset Matching with Refined Covariate Balance***Description**

Tools for optimal subset matching of treated units and control units in observational studies, with support for refined covariate balance constraints, (including fine and near-fine balance as special cases). A close relative is the 'rcbalance' package. See Pimentel, et al.(2015) <doi:10.1080/01621459.2014.997879> and Pimentel and Kelz (2020) <doi:10.1080/01621459.2020.1720693>.

Details

The DESCRIPTION file:

```
Package:      rcbsubset
Type:         Package
Title:        Optimal Subset Matching with Refined Covariate Balance
Version:      1.1.5
Date:         2021-10-15
Author:       Samuel D. Pimentel
Maintainer:   Samuel D. Pimentel <spi@berkeley.edu>
Description:  Tools for optimal subset matching of treated units and control units in observational studies, with support for re
Depends:      R (>= 3.2.0), MASS, plyr
Imports:      rcbalance
License:      MIT + file LICENSE
Suggests:    optmatch, testthat
```

Index of help topics:

```
callrelax      Solve Network Flow Problem using RELAX-IV
                Algorithm
dist2net       Building and Manipulating Network Flow Problems
rcbsubset      Optimal Matching with Refined Covariate Balance
rcbsubset-package Optimal Subset Matching with Refined Covariate
                Balance
```

This package computes 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. In addition, treated units may be excluded in an optimal manner (using a penalty parameter) to improve the quality of the match. For more detail see the references.

The main function is `rcbsubset`, which takes a distance/sparsity matrix or matrix-like object containing information about matchability of the treated and control units and a list of fine balance variables and produces a match. The other functions are largely for internal use and should not be

needed by the large majority of users. The syntax and code structure is very similar in the closely related antecedent package `rcbalance`, which provides more helper functions for constructing matches but does not support optimal subset matching.

IMPORTANT NOTE: the functionality of this package is greatly reduced unless the `optmatch` package ($v \geq 0.9-1$) is also loaded. When attempting to run the `rcbalance` command without having loaded `optmatch`, the users will receive an error message. The first reference below gives background on `optmatch`.

Author(s)

Samuel D. Pimentel

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

References

Hansen, B.B. and Klopfer, S.O. (2006) Optimal full matching and related designs via network flows, *JCGS* 15 609-627.

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., and Kelz, R.R. (2020). Optimal tradeoffs in matched designs comparing US-trained and internationally trained surgeons. *JASA* 115 (532), 1675-1688.

callrelax

Solve Network Flow Problem using RELAX-IV Algorithm

Description

An `rcbalance` method not meant to be called directly by users. Solves network flow optimization problems by calling the RELAX-IV algorithm, as implemented in FORTRAN by Dimitri Bertsekas and Paul Tseng.

IMPORTANT NOTE 1: the RELAX-IV code is not contained in this R package due to software license issues. Users can only access it by loading the `optmatch` package ($\geq 0.9-1$) and accepting its license. The reference below gives background on `optmatch`.

Usage

```
callrelax(net)
```

Arguments

`net` 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):

- `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.

Value

A list with the following elements:

<code>crash</code>	an integer, equal to zero if the algorithm ran correctly and equal to 1 if it crashed.
<code>feasible</code>	an integer, equal to zero if the problem is infeasible and equal to 1 if it is feasible.
<code>x</code>	a vector equal in length to the number of arcs in argument problem <code>net</code> , 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

References

Hansen, B.B. and Klopfer, S.O. (2006) Optimal full matching and related designs via network flows, JCGS 15 609-627.

dist2net

Building and Manipulating Network Flow Problems

Description

These are internal `rcsubset` 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 `rcsubset` method.

Usage

```
dist2net(dist.struct, k, exclude.treated = FALSE, exclude.penalty = NULL, ncontrol = NULL,
tol = 1e-3)
```

```
dist2net.matrix(dist.struct, k, exclude.treated = FALSE, exclude.penalty = NULL,
tol = 1e-3)
```

```
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>exclude.penalty</code>	a parameter that gives the cost of excluding a treated unit. If left NULL it will be set to a very large value designed to ensure treated units are never excluded if they can be matched. Lower values may result in subsets of treated units being excluded.
<code>ncontrol</code>	the total number of controls in the matching problem. If left NULL the function will attempt to compute it by counting controls referenced in the distance object provided.
<code>tol</code>	edge cost tolerance. This is the smallest tolerated difference between matching costs. It is used in these functions only when <code>exclude.penalty</code> is NULL, to choose a default penalty that is small enough not to cause integer overflows.
<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 <code>value</code> 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 <code>value</code> 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 `nnodes` 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

rcbsubset

Optimal Matching with Refined Covariate Balance

Description

This function computes an optimal match with refined covariate balance.

Usage

```
rcbsubset(distance.structure, near.exact = NULL, fb.list = NULL,  
treated.info = NULL, control.info = NULL, exclude.penalty = NULL, penalty = 2, tol = 1e-3)
```

Arguments

- `distance.structure` A distance matrix, a sparse distance matrix of class `InfinitySparseMatrix` (produced by various functions in the `optmatch` package), or a a list of vectors that encodes information about covariate distances between treated and control units (produced by the `build.dist.struct` function in the `rcbalance` package). If a matrix is given, rows should correspond to treated units and columns to control units. Distances of `Inf` should be used to indicate that the units in question must never be matched to each other.
- `near.exact` 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.
- `fb.list` 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 `fb.list` argument is specified, the `treated.info` and `control.info` arguments must also be specified.
- `treated.info` 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 `distance.structure` argument, and it is assumed that row `i` contains covariate information for the treated unit described by element `i` of `distance.structure`. In addition, the column count and column names must be identical to those of the `control.info` argument, and the column names must include all of the covariate names mentioned in the `near.exact` and `fb.list` arguments.
- `control.info` 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 `distance.structure` argument, and it is assumed that row `i` contains the covariate information for the control indexed by `i` in `distance.structure`. In addition, the column count and column names must be identical to those of the `treated.info` argument.
- `exclude.penalty` A parameter that gives the cost of excluding a treated unit. If left `NULL` it will be set to a very large value designed to ensure treated units are never excluded if

	they can be matched. Lower values may result in subsets of treated units being excluded.
penalty	a nonnegative value. 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. It is not recommended that the user change this parameter from its default value.
tol	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.

Details

In order to perform matching, `rcbalance` requires the user to load the `optmatch` ($\geq 0.9-1$) package separately. The manual loading is required due to software license issues. If the package is not loaded the `rcbalance` command will fail with an error saying the `optmatch` package is not present.

The second reference below gives background on `optmatch`.

Value

A list with the following components:

matches	a nt by k matrix containing the matched sets produced by the algorithm (where nt 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.
fb.tables	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

Hansen, B.B. and Klopfer, S.O. (2006) Optimal full matching and related designs via network flows, *JCGS* 15 609-627.

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.

Rosenbaum, P.R. (2012) Optimal matching of an optimally chosen subset in observational studies, *JCGS* 21.1: 57-71.

Examples

```
library(optmatch)
data(nuclearplants)

#match exactly on 3 binaries
exact.mask <- exactMatch(pr ~ pt + ct + bw, data = nuclearplants)
my.dist.matrix <- match_on(pr ~ date + t1 + t2 + cap + ne + cum.n,
within = exact.mask, data = nuclearplants)

#one treated unit out of 10 is excluded
rcbsubset(my.dist.matrix)

#repeat under a refined balance constraint
rcbsubset(my.dist.matrix, fb.list = list('ne'),
treated.info = nuclearplants[which(nuclearplants$pr ==1),],
control.info = nuclearplants[which(nuclearplants$pr == 0),])

#specifying a low exclude.penalty leads to more individuals excluded
rcbsubset(my.dist.matrix, fb.list = list('ne'),
treated.info = nuclearplants[which(nuclearplants$pr ==1),],
control.info = nuclearplants[which(nuclearplants$pr == 0),], exclude.penalty = 1.5)

## Not run:
#match using distance objects created by rcbalance package
library(rcbalance)

my.dist.struct <- build.dist.struct(z = nuclearplants$pr,
X = subset(nuclearplants[c('date', 't1', 't2', 'cap', 'ne', 'cum.n')]),
exact = paste(nuclearplants$pt, nuclearplants$ct, nuclearplants$bw, sep = '.'))

rcbsubset(my.dist.struct, fb.list = list('ne'),
treated.info = nuclearplants[which(nuclearplants$pr ==1),],
control.info = nuclearplants[which(nuclearplants$pr == 0),], exclude.penalty = 15)

## End(Not run)
```

Index

`add.layer (dist2net)`, 4

`callrelax`, 3

`dist2net`, 4

`penalize.near.exact (dist2net)`, 4

`penalty.update (dist2net)`, 4

`rcsubset`, 6

`rcsubset-package`, 2

`remove.layer (dist2net)`, 4