

Package ‘rcbalance’

September 2, 2014

Type Package

Title Large, sparse optimal matching with refined covariate balance

Version 1.0

Date 2014-08-08

Author Samuel D. Pimentel

Maintainer Samuel D. Pimentel <spi@wharton.upenn.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.

Depends R (>= 3.0.0), MASS, plyr

License MIT + file LICENSE

NeedsCompilation no

Repository CRAN

Date/Publication 2014-09-02 22:18:56

R topics documented:

rcbalance-package	2
build.dist.struct	2
callrelax	3
dist2net	5
rcbalance	6

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

rcbalance-package *Sparse Optimal Matching with Refined Covariate Balance*

Description

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 reference.

Details

Package: rcbalance
 Type: Package
 Version: 1.0
 Date: 2014-07-10
 License: What license is it under?

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 other functions are largely for internal use and should not be needed by the large majority of users.

IMPORTANT NOTE: the functionality of this package is greatly reduced if 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 reference below gives background on `optmatch`.

Author(s)

Samuel D. Pimentel <spi@wharton.upenn.edu>

References

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

build.dist.struct *Build Distance Structure for Matching with Refined Balance*

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)
```

Arguments

z	a vector of treatment and control indicators, 1 for treatment and 0 for control.
X	a data frame 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'.

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](#)

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	<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.
-----	--

Value

A list with the following elements:

crash	an integer, equal to zero if the algorithm ran correctly and equal to 1 if it crashed.
feasible	an integer, equal to zero if the problem is infeasible and equal to 1 if it is feasible.
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

References

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

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)
```

```
dist2net.matrix(dist.struct, k)
```

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

```
remove.layer(net.layers, layer.idx)
```

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

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>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>layer.idx</code>	number of the layer to be removed from the network. Layers are numbered from coarsest to finest.
<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).

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) and `remove.layer` removes these layers. `penalty.update` is used to change the penalties for each layer. 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, fb.list = NULL,
          treated.info = NULL, control.info = NULL, k = 1,
          penalty = 3)
```

Arguments

- `distance.structure` 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 `nc` control units in the problem and that they are numbered from 1 to `nc`. The names of each vector in the list give the index (in the vector `1:nc`) 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 `nc` and rownames 1 through `nc`.
- Alternatively, this same information can be passed as a matrix or `InfinitySparseMatrix` with rows corresponding to treated units and columns corresponding to controls. Entries given as `Inf` correspond to pairs that cannot be matched. When working with very large datasets, however, it is recommended to use the list-of-vectors distance specification.
- Note that non-integer distances in `distance.structure` will be rounded to the nearest integer, so users may wish to multiply their distance by a factor of 100 or 1000 to preserve fine distinctions.
- `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 `fb.list` argument.
- `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.

k	a nonnegative integer. The number of control units to which each treated unit will be matched.
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.

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 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.

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)`, 5

`build.dist.struct`, 2

`callrelax`, 3

`dist2net`, 5

`penalty.update (dist2net)`, 5

`rcbalance`, 3, 6

`rcbalance-package`, 2

`remove.layer (dist2net)`, 5