

# Package ‘nbpMatching’

April 19, 2016

**Type** Package

**Title** Functions for Optimal Non-Bipartite Matching

**Version** 1.5.1

**Date** 2016-04-18

**Author** Cole Beck [aut, cre], Bo Lu [aut], Robert Greevy [aut]

**Maintainer** Cole Beck <cole.beck@vanderbilt.edu>

**Description** Perform non-bipartite matching and matched randomization. A ``bipartite" matching utilizes two separate groups, e.g. smokers being matched to nonsmokers or cases being matched to controls. A ``non-bipartite" matching creates mates from one big group, e.g. 100 hospitals being randomized for a two-arm cluster randomized trial or 5000 children who have been exposed to various levels of secondhand smoke and are being paired to form a greater exposure vs. lesser exposure comparison. At the core of a non-bipartite matching is a  $N \times N$  distance matrix for  $N$  potential mates. The distance between two units expresses a measure of similarity or quality as mates (the lower the better). The 'gendistance()' and 'distancematrix()' functions assist in creating this. The 'nonbimatch()' function creates the matching that minimizes the total sum of distances between mates; hence, it is referred to as an ``optimal" matching. The 'assign.grp()' function aids in performing a matched randomization. Note bipartite matching can be performed using the prevent option in 'gendistance()'.

**License** GPL (>= 2)

**Imports** methods, stats, Hmisc, MASS, utils

**URL** <http://biostat.mc.vanderbilt.edu/wiki/Main/MatchedRandomization>

**Collate** 'distancematrix.R' 'nonbimatch.R' 'assign.grp.R'  
'fill.missing.R' 'gendistance.R' 'make.phantoms.R'  
'nbpmatching-package.R' 'qom.R' 'quantile.R' 'subsetMatches.R'  
'utils.R'

**RoxygenNote** 5.0.1

**NeedsCompilation** yes

**Repository** CRAN

**Date/Publication** 2016-04-19 08:12:05

**R topics documented:**

nbpMatching-package . . . . .	2
assign.grp . . . . .	3
distancematrix . . . . .	4
fill.missing . . . . .	6
gendistance . . . . .	7
get.sets . . . . .	9
make.phantoms . . . . .	10
nonbimatch . . . . .	11
qom . . . . .	12
quantile,distancematrix-method . . . . .	14
scalar.dist . . . . .	15
subsetMatches . . . . .	15

<b>Index</b>	<b>17</b>
--------------	-----------

---

nbpMatching-package    *Nonbipartite Matching*

---

**Description**

This package will take an input distance matrix and generate the set of pairwise matches that minimizes the sum of distances between the pairs by running nonbimatch.

**Details**

The most current documentation is available at <http://biostat.mc.vanderbilt.edu/wiki/Main/MatchedRandomization>.

**Author(s)**

Bo Lu, Robert Greevy, Cole Beck

Maintainer: Cole Beck <cole.beck@vanderbilt.edu>

**References**

Lu B, Greevy R, Xu X, Beck C. Optimal Nonbipartite Matching and its Statistical Applications. The American Statistician. Vol. 65, no. 1. : 21-30. 2011.

Greevy RA Jr, Grijalva CG, Roumie CL, Beck C, Hung AM, Murff HJ, Liu X, Griffin MR. Reweighted Mahalanobis distance matching for cluster-randomized trials with missing data. Pharmacoepidemiol Drug Saf. 2012 May;21 Suppl 2:148-54. doi: 10.1002/pds.3260.

## Examples

```
# create a covariate matrix
df <- data.frame(id=LETTERS[1:25], val1=rnorm(25), val2=rnorm(25))
# create distances
df.dist <- gendistance(df, idcol=1)
# create distancematrix object
df.mdm <- distancematrix(df.dist)
# create matches
df.match <- nonbimatch(df.mdm)
# review quality of matches
df.qom <- qom(df.dist$cov, df.match$matches)

# some helper functions are available
# runner -- start with the covariate, run through the entire process
df.1 <- runner(df, idcol=1)
# full.qom -- start with the covariate, generate a full quality of match report
df.2 <- full.qom(df)

## Not run:
try a large matrix
nonbimatch(distancematrix(as.matrix(dist(sample(1:10^8, 5000, replace=TRUE))))))

## End(Not run)
```

---

assign.grp

*Random Group Assignment*

---

## Description

Randomly assign each element into treatment group A or B.

## Usage

```
assign.grp(matches, seed = 68, ...)
```

## Arguments

matches	A data.frame or nonbimatch object. Contains information on how to match the covariate data set.
seed	Seed provided for random-number generation. Default value of 68.
...	Additional arguments, not used at the moment.

## Details

This function takes the matched pairs generated by nonbimatch and randomly assigns each element to a group.

**Value**

original data.frame with treatment group column

**Author(s)**

Cole Beck

**See Also**

[nonbimatch](#)

**Examples**

```
df <- data.frame(id=LETTERS[1:25], val1=rnorm(25), val2=rnorm(25))
df.dist <- gendistance(df, idcol=1)
df.mdm <- distancematrix(df.dist)
df.match <- nonbimatch(df.mdm)
assign.grp(df.match)
assign.grp(df.match$matches)
```

---

distancematrix

*Distance matrix*

---

**Description**

The `distancematrix` function is used to reformat the input distance matrix into the format required by the nonbipartite matching Fortran code. The original matrix should have dimensions  $N \times N$ , where  $N$  is the total number of elements to be matched. The matrix may be created in R and input into the `distancematrix` function. Alternately, the matrix may be read in from a CSV file, i.e. a text file where distances in a given row are delimited by commas. If a list element is given, it should have a data.frame element named "dist", preferably generated by the `gendistance` function.

**Usage**

```
distancematrix(x, ...)
```

**Arguments**

**x** A matrix, data.frame, list or filename. This should be an  $N \times N$  distance matrix for the  $N$  elements to be matched. The values in the diagonal are ignored because an element cannot be matched to itself. Using zeros in the diagonal is preferable, although other values are acceptable provided they are not so large that they distort the scaling of the other values.

**...** Additional arguments, potentially used when reading in a filename and passed into `read.csv`.

## Details

- The `distancematrix` function is used to reformat the input distance matrix into the format required by the nonbipartite matching Fortran code.
- If an extra column or row is present, it will be converted into row names. In other words, if the matrix has dimensions  $(N + 1) \times N$ , or  $N \times (N + 1)$ , then the function will take the first row, or column, as an ID column. If both row and column names are present, i.e. a  $(N + 1) \times (N + 1)$  matrix, the function cannot identify the names.
- If an odd number of elements exist, a ghost element, or sink, will be created whose distance is zero to all of the other elements. For example, when matching 17 elements, the function will create an 18th element that matches every element perfectly. This sink may or not be appropriate for your application. Naturally, you may create sinks as needed in the distance matrix you input to the `distancematrix` function.
- The elements of `distancematrix` may not be re-assigned once created. In other words, you cannot edit the formatted distance matrix. You need to edit the matrix being input into the `distancematrix` function.

## Value

`distancematrix` S4 object

## Author(s)

Cole Beck

## See Also

[nonbimatch](#) [gendistance](#)

## Examples

```
plainmatrix<-as.matrix(dist(sample(1:25, 8, replace=TRUE)))
diag(plainmatrix) <- 99999 # setting diagonal to an infinite distance for
                          # pedagogical reasons (the diagonal may be left
                          # as zero)
mdm<-distancematrix(plainmatrix)
df <- data.frame(id=LETTERS[1:25], val1=rnorm(25), val2=rnorm(25))
df[sample(seq_len(nrow(df)), ceiling(nrow(df)*0.1)), 2] <- NA
df.dist <- gendistance(df, idcol=1, ndiscard=2)
mdm2 <- distancematrix(df.dist)
```

---

`fill.missing`*Data Imputation*

---

### Description

The `fill.missing` function uses the `transcan` function from the **Hmisc** package to impute values for the given `data.frame`.

### Usage

```
fill.missing(x, seed = 101, simplify = TRUE, idcol = "id", ...)
```

### Arguments

<code>x</code>	A <code>data.frame</code> object. It should have missing values.
<code>seed</code>	Seed provided for random-number generation. Default value of 101.
<code>simplify</code>	logical: whether to remove duplicate missingness columns.
<code>idcol</code>	An integer value or character string. Indicates the column containing IDs, specified as column index or column name. Defaults to "id", or NA, when not found.
<code>...</code>	Additional arguments, potentially passed to <code>transcan</code> .

### Details

The `fill.missing` function will fill the missing values within a `data.frame` with the values imputed with the `transcan` function. An `idcol` may be specified to prevent including the use of IDs in the imputation. In addition for every column that contains missing data, a new column will be attached to the `data.frame` containing an indicator of missingness. A "1" indicates that the value was missing and has been imputed.

### Value

`data.frame` with imputed values

### Author(s)

Cole Beck

### See Also

[transcan](#)

**Examples**

```
set.seed(1)
df <- data.frame(id=LETTERS[1:25], val1=rnorm(25), val2=rnorm(25))
df[sample(seq_len(nrow(df)), ceiling(nrow(df)*0.1)), 2] <- NA
df <- fill.missing(df)
```

---

gendistance	<i>Generate a Distance Matrix</i>
-------------	-----------------------------------

---

**Description**

The `gendistance` function creates an  $(N + K) \times (N + K)$  distance matrix from an  $N \times P$  covariates matrix, where  $N$  is the number of subjects,  $P$  the number of covariates, and  $K$  the number of phantom subjects requested (see `ndiscard` option). Provided the covariates' covariance matrix is invertible, the distances computed are Mahalanobis distances, or if covariate weights are provided, Reweighted Mahalanobis distances (see `weights` option and Greevy, et al., *Pharmacoepidemiology and Drug Safety* 2012).

**Usage**

```
gendistance(covariate, idcol = NULL, weights = NULL, prevent = NULL,
  force = NULL, rankcols = NULL, missing.weight = 0.1, ndiscard = 0,
  singular.method = "solve", talisman = NULL, prevent.res.match = NULL,
  ...)
```

**Arguments**

<code>covariate</code>	A <code>data.frame</code> object, containing the covariates of the data set.
<code>idcol</code>	An integer or column name, providing the index of the column containing row ID's.
<code>weights</code>	A numeric vector, the length should match the number of columns. This value determines how much weight is given to each column when generating the distance matrix.
<code>prevent</code>	A vector of integers or column names, providing the index of columns that should be used to prevent matches. When generating the distance matrix, elements that match on these columns are given a maximum distance.
<code>force</code>	An integer or column name, providing the index of the column containing information used to force pairs to match.
<code>rankcols</code>	A vector of integers or column names, providing the index of columns that should have the rank function applied to them before generating the distance matrix.

<code>missing.weight</code>	A numeric value, or vector, used to generate the weight of missingness indicator columns. Missingness indicator columns are created if there is missing data within the data set. Defaults to 0.1. If a single value is supplied, weights are generating by multiplying this by the original columns' weight. If a vector is supplied, it's length should match the number of columns with missing data, and the weight is used as is.
<code>ndiscard</code>	An integer, providing the number of elements that should be allowed to match phantom values. The default value is 0.
<code>singular.method</code>	A character string, indicating the function to use when encountering a singular matrix. By default, <code>solve</code> is called. The alternative is to call <code>ginv</code> from the <b>MASS</b> package.
<code>talisman</code>	An integer or column name, providing location of talisman column. The talisman column should only contains values of 0 and 1. Records with zero will match phantoms perfectly, while other records will match phantoms at max distance.
<code>prevent.res.match</code>	An integer or column name, providing location of the column containing assigned treatment groups. This is useful in some settings, such as trickle-in randomized trials. When set, non-NA values from this column are replaced with the value 1. This prevents records with previously assigned treatments (the 'reservior') from matching each other.
<code>...</code>	Additional arguments, not used at this time.

## Details

Given a `data.frame` of covariates, generate a distance matrix. Missing values are imputed with `fill.missing`. For each column with missing data, a missingness indicator column will be added. Phantoms are fake elements that perfectly match all elements. They can be used to discard a certain number of elements.

## Value

a list object with several elements

<code>dist</code>	generated distance matrix
<code>cov</code>	covariate matrix used to generate distances
<code>ignored</code>	ignored columns from original covariate matrix
<code>weights</code>	weights applied to each column in covariate matrix
<code>prevent</code>	columns used to prevent matches
<code>mates</code>	index of rows that should be forced to match
<code>rankcols</code>	index of columns that should use rank
<code>missing.weight</code>	weight to apply to missingness indicator columns
<code>ndiscard</code>	number of elements that will match phantoms



**Author(s)**

Cole Beck

**See Also**[distancematrix](#)**Examples**

```
set.seed(1)
df <- data.frame(id=LETTERS[1:25], val1=rnorm(25), val2=rnorm(25))
# add some missing data
df[sample(seq_len(nrow(df)), ceiling(nrow(df)*0.1)), 2] <- NA
df.dist <- gendistance(df, idcol=1, ndiscard=2)
# up-weight the second column
df.weighted <- gendistance(df, idcol=1, weights=c(1,2,1), ndiscard=2, missing.weight=0.25)
df[,3] <- df[,2]*2
df.solve <- gendistance(df, idcol=1, ndiscard=2)
df.solve.ginv <- gendistance(df, idcol=1, ndiscard=2, singular.method="ginv")
```

---

`get.sets`*Get named sets of matches*

---

**Description**

Create a factor variable using the names from a matched data set.

**Usage**

```
get.sets(matches, remove.unpaired = TRUE, ...)
```

**Arguments**

<code>matches</code>	A data.frame or nonbimatch object. Contains information on how to match the covariate data set.
<code>remove.unpaired</code>	A boolean value. The default is to remove elements matched to phantom elements.
<code>...</code>	Additional arguments, not used at this time.

**Details**

Calculate a name for each pair by using the ID columns from the matched data set. Return a factor of these named pairs.

**Value**

a factor vector

**Author(s)**

Jake Bowers, <http://jakebowers.org/>, Cole Beck

**Examples**

```
df <- data.frame(id=LETTERS[1:25], val1=rnorm(25), val2=rnorm(25))
df.dist <- gendistance(df, idcol=1)
df.mdm <- distancematrix(df.dist)
df.match <- nonbimatch(df.mdm)
get.sets(df.match)
get.sets(df.match$matches)
# include the phantom match
get.sets(df.match$matches, FALSE)
```

---

make.phantoms

*Add Phantom Rows and Columns*

---

**Description**

The `make.phantoms` function will take an  $N \times N$  matrix and add  $NP$  phantom elements, thus creating a matrix with  $N + NP \times N + NP$  dimensions.

**Usage**

```
make.phantoms(x, nphantoms, name = "phantom", maxval = Inf, ...)
```

**Arguments**

<code>x</code>	A matrix or data.frame object, with $N \times N$ dimensions.
<code>nphantoms</code>	An integer, providing the number of phantom elements to add.
<code>name</code>	A character string, indicating the name attribute for new elements. Defaults to "phantom".
<code>maxval</code>	An integer value, the default value to give the pairs of phantoms (indices $[N+1:N+NP, N+1:N+NP]$ ), assumed to be a maximum distance. Defaults to Inf.
<code>...</code>	Additional arguments, not used at this time.

**Details**

This function is internal to the `gendistance` function, but may be useful in manufacturing personalized distance matrices. Phantoms are fake elements that perfectly match all elements. They can be used to discard a certain number of elements.

**Value**

a matrix or data.frame object

**Author(s)**

Cole Beck

**See Also**

[gendistance](#) [distancematrix](#)

**Examples**

```
# 5x5 distance matrix
dist.mat <- matrix(c(0,5,10,15,20,5,0,15,25,35,10,15,0,25,40,15,25,25,0,15,20,35,40,15,0), nrow=5)
# add one phantom element
dm.ph <- make.phantoms(dist.mat, 1)
# create distancematrix object
distancematrix(dm.ph)
# add three phantoms
make.phantoms(dist.mat, 3)
```

---

nonbimatch

*Nonbipartite Matching*


---

**Description**

The nonbinmatch function creates the set of pairwise matches that minimizes the sum of distances between the pairs.

**Usage**

```
nonbimatch(mdm, threshold = NA, precision = 6, ...)
```

**Arguments**

mdm	A distancematrix object. See the distancematrix function.
threshold	An numeric value, indicating the distance needed to create chameleon matches.
precision	The largest value in the matrix will have at most this many digits. The default value is six.
...	Additional arguments, these are not used.

**Details**

The nonbinmatch function calls the Fortran code (Derigs) and set of pairwise matches that minimizes the sum of distances between the pairs.

**Value**

nonbimatch S4 object with several elements

matches	data.frame containing matches
halves	data.frame containing each match
total	sum of the distances across all pairs
mean	mean distance for each pair

**Author(s)**

Cole Beck

**See Also**

[distancematrix](#)

**Examples**

```
plainmatrix<-as.matrix(dist(sample(1:25, 8, replace=TRUE)))
diag(plainmatrix) <- 99999 # setting diagonal to an infinite distance for
                           # pedagogical reasons (the diagonal may be left
                           # as zero)
mdm<-distancematrix(plainmatrix)
res<-nonbimatch(mdm)
```

---

qom

*Quality of Match*

---

**Description**

Quality of matches show how well matched pairs differ. For each variable the average distance is generated. Each item in a pair is assigned a group and after several iterations the quantile of these average distances is returned.

**Usage**

```
qom(covariate, matches, iterations = 10000, probs = NA, use.se = FALSE,
     all.vals = FALSE, seed = 101, ...)
```

**Arguments**

covariate	A data.frame object.
matches	A data.frame or nonbimatch object. Contains information on how to match the covariate data set.
iterations	An integer. Number of iterations to run, defaults to 10,000.
probs	A numeric vector. Probabilities to pass to the quantile function.
use.se	A logical value. Determines if the standard error should be computed. Default value of FALSE.
all.vals	A logical value. Determines if false matches should be included in comparison. Default value of FALSE.
seed	Seed provided for random-number generation. Default value of 101.
...	Additional arguments, not used at the moment.

**Details**

This function is useful for determining the effectiveness of your weights (when generating a distance matrix). Weighting a variable more will lower the average distance, but it could penalize the distance of the other variables. Calculating the standard error requires calling `hdquantile` from **Hmisc**. The quantiles may be slightly different when using `hdquantile`.

**Value**

a list object containing elements with quality of match information

q	data.frame with quantiles for each covariate
se	data.frame with standard error for each covariate
sd	vector with standard deviate for each covariate

**Author(s)**

Cole Beck

**Examples**

```
df <- data.frame(id=LETTERS[1:25], val1=rnorm(25), val2=rnorm(25))
df.dist <- gendistance(df, idcol=1)
df.mdm <- distancematrix(df.dist)
df.match <- nonbimatch(df.mdm)
qom(df.dist$cov, df.match)
qom(df.dist$cov, df.match$matches)
```

---

quantile,distancematrix-method

*Quantile for upper-triangular values in distance matrix*

---

### Description

Extend the stats `quantile` function for handling `distancematrix` objects.

### Usage

```
## S4 method for signature 'distancematrix'  
quantile(x, probs, ...)
```

### Arguments

<code>x</code>	A <code>distancematrix</code> object.
<code>probs</code>	numeric vector or probabilities with values in [0,1].
<code>...</code>	Additional arguments, passed to <code>quantile</code> .

### Details

The upper.triangular values of the distance matrix object are passed to the `quantile` function.

### Value

numeric vector of quantiles corresponding to the given probabilities

### Author(s)

Cole Beck

### Examples

```
plainmatrix<-as.matrix(dist(sample(1:25, 8, replace=TRUE)))  
mdm<-distancematrix(plainmatrix)  
quantile(mdm, probs=c(0.0, 0.25, 0.50, 0.75, 1.00))
```

---

scalar.dist	<i>Calculate scalar distance</i>
-------------	----------------------------------

---

**Description**

Calculate the scalar distance between elements of a matrix.

**Usage**

```
scalar.dist(x, ...)
```

**Arguments**

x	A vector of numeric values.
...	Additional arguments, not used at this time.

**Details**

Take the absolute difference between all elements in a vector, and return a matrix of the distances.

**Value**

a matrix object

**Author(s)**

Jake Bowers, <http://jakebowers.org/>, Cole Beck

**Examples**

```
scalar.dist(1:10)
```

---

subsetMatches	<i>Subset Matches</i>
---------------	-----------------------

---

**Description**

Remove unpaired or unnecessary matches.

**Usage**

```
subsetMatches(matches, phantom = TRUE, chameleon = TRUE, ghost = TRUE,  
infinite = TRUE, halvesOnly = TRUE)
```

**Arguments**

matches	A nonbimatch object.
phantom	A logical value. Remove elements matched to phantom elements.
chameleon	A logical value. Remove elements matched to chameleon elements.
ghost	A logical value. Remove elements matched to ghost elements.
infinite	A logical value. Remove elements matched at infinite distance. This will include elements forced to match in spite of having an infinite distance set by the prevent option in <a href="#">gendistance</a> .
halvesOnly	A logical value. Use halves element instead of matches.

**Details**

Given a nonbimatch object, remove elements matched to phantoms, chameleons, or ghosts. Also remove pairs whose distance is infinite.

**Value**

a data.frame

**Author(s)**

Cole Beck

**Examples**

```
df <- data.frame(id=LETTERS[1:25], val1=rnorm(25), val2=rnorm(25))
df.dist <- gendistance(df, idcol=1, ndiscard=4)
df.mdm <- distancematrix(df.dist)
df.match <- nonbimatch(df.mdm)
subsetMatches(df.match)
subsetMatches(df.match, halvesOnly=FALSE)
subsetMatches(df.match, phantom=FALSE)
```



# Index

- \*Topic **array**
  - nbpMatching-package, 2
- \*Topic **cluster**
  - nbpMatching-package, 2
- \*Topic **package**
  - nbpMatching-package, 2
- assign.grp, 3
- assign.grp,data.frame-method
  - (assign.grp), 3
- assign.grp,nonbimatch-method
  - (assign.grp), 3
- distancematrix, 4, 9, 11, 12
- distancematrix,character-method
  - (distancematrix), 4
- distancematrix,data.frame-method
  - (distancematrix), 4
- distancematrix,list-method
  - (distancematrix), 4
- distancematrix,matrix-method
  - (distancematrix), 4
- distancematrix-class (distancematrix), 4
- fill.missing, 6, 8
- fill.missing,data.frame-method
  - (fill.missing), 6
- gendistance, 5, 7, 10, 11, 16
- gendistance,data.frame-method
  - (gendistance), 7
- get.sets, 9
- get.sets,data.frame-method (get.sets), 9
- get.sets,nonbimatch-method (get.sets), 9
- ginv, 8
- hdquantile, 13
- make.phantoms, 10
- make.phantoms,ANY,missing-method
  - (make.phantoms), 10
- make.phantoms,data.frame,numeric-method
  - (make.phantoms), 10
- make.phantoms,matrix,numeric-method
  - (make.phantoms), 10
- nbpMatching (nbpMatching-package), 2
- nbpMatching-package, 2
- nonbimatch, 4, 5, 11
- nonbimatch,distancematrix-method
  - (nonbimatch), 11
- nonbimatch-class (nonbimatch), 11
- qom, 12
- qom,data.frame,data.frame-method (qom), 12
- qom,data.frame,nonbimatch-method (qom), 12
- quantile, 14
- quantile
  - (quantile,distancematrix-method), 14
- quantile,distancematrix-method, 14
- scalar.dist, 15
- scalar.dist,vector-method
  - (scalar.dist), 15
- solve, 8
- subsetMatches, 15
- subsetMatches,nonbimatch-method
  - (subsetMatches), 15
- transcan, 6