

Package ‘reclin2’

January 7, 2022

Type Package

Title Record Linkage Toolkit

Version 0.1.1

Description Functions to assist in performing probabilistic record linkage and deduplication: generating pairs, comparing records, em-algorithm for estimating m- and u-probabilities (I. Fellegi & A. Sunter (1969) <[doi:10.1080/01621459.1969.10501049](https://doi.org/10.1080/01621459.1969.10501049)>, T.N. Herzog, F.J. Scheuren, & W.E. Winkler (2007), ``Data Quality and Record Linkage Techniques'', ISBN:978-0-387-69502-0), forcing one-to-one matching. Can also be used for pre- and post-processing for machine learning methods for record linkage. Focus is on memory, CPU performance and flexibility.

BugReports <https://github.com/djvanderlaan/reclin2/issues>

URL <https://github.com/djvanderlaan/reclin2>

Depends data.table, R (>= 3.6.0)

Imports stringdist, stats, utils, lpSolve, Rcpp, parallel

Suggests simplermarkdown

LinkingTo Rcpp

VignetteBuilder simplermarkdown

SystemRequirements C++11

License GPL-3

LazyLoad yes

Encoding UTF-8

RoxygenNote 7.1.1

NeedsCompilation yes

Author Jan van der Laan [aut, cre] (<<https://orcid.org/0000-0002-0693-1514>>)

Maintainer Jan van der Laan <r@eoos.dds.nl>

Repository CRAN

Date/Publication 2022-01-07 13:22:42 UTC

R topics documented:

add_from_x	2
cluster_call	3
cluster_collect	4
cluster_modify_pairs	5
cluster_pair	6
cluster_pair_blocking	7
cluster_pair_minsim	9
compare_pairs.cluster_pairs	10
compare_vars.cluster_pairs	12
deduplicate_equivalence	13
get_inspect_pairs	14
greedy	15
identical	15
link	17
linkexample1	18
match_n_to_m	18
pair	19
pair_blocking	20
pair_minsim	21
predict.problink_em	22
problink_em	24
select_greedy.cluster_pairs	25
select_threshold.cluster_pairs	28
summary.problink_em	29
tabulate_patterns.cluster_pairs	30
town_names	31
Index	32

add_from_x	<i>Add a variable from one of the data sets to pairs</i>
------------	----------------------------------------------------------

Description

Add a variable from one of the data sets to pairs

Usage

```
add_from_x(pairs, variable, new_variable = variable, ...)
```

```
add_from_y(pairs, variable, new_variable = variable, ...)
```

Arguments

pairs	data.table with pairs. Should contain the columns .x and .y.
variable	name of the variable that should be added
new_variable	optional variable name of the new variable in pairs. When omitted variable is used.
...	other parameters are passed on to compare_vars. Especially inplace, x and y might be of interest.

Value

Returns the pairs with the column added. When inplace = TRUE pairs is returned invisibly and the original pairs is modified.

cluster_call	<i>Call a function on each of the worker nodes and pass it the pairs</i>
--------------	--------------------------------------------------------------------------

Description

Call a function on each of the worker nodes and pass it the pairs

Usage

```
cluster_call(pairs, fun, ...)
```

Arguments

pairs	an object or type cluster_pairs as created for example by cluster_pair .
fun	a function to call on each of the worker nodes. See details on the arguments of this function.
...	additional arguments are passed on to fun.

Details

The function will have to accept the following arguments as its first three arguments:

pairs the data.table with the pairs of the worker node.

x a data.table with the portion of x present on the worker node.

y a data.table with y.

Value

The function will return a list with for each worker the result of the function call. When the functions return NULL the result is returned invisibly. Because the result is returned to main node, make sure you don't accidentally return all pairs. If you don't want to return anything end your function with NULL.

Examples

```

# Generate some pairs
library(parallel)
data("linkexample1", "linkexample2")
cl <- makeCluster(2)
pairs <- cluster_pair(cl, linkexample1, linkexample2)
compare_pairs(pairs, c("lastname", "firstname", "address", "sex"))

# Add a new column to pairs
cluster_call(pairs, function(pairs, ...) {
  pairs[, name := firstname & lastname]
  # we don't want to return the pairs; so make sure to return something
  # else
  NULL
})

# Get the number of pairs on each node
lengths <- cluster_call(pairs, function(pairs, ...) {
  nrow(pairs)
})
lengths <- unlist(lengths)
lengths

# Cleanup
stopCluster(cl)

```

cluster_collect	<i>Collect pairs from cluster nodes</i>
-----------------	-----------------------------------------

Description

Collect pairs from cluster nodes

Usage

```
cluster_collect(pairs, select = NULL, clear = FALSE)
```

Arguments

pairs	an object or type cluster_pairs as created for example by cluster_pair .
select	the name of a logical column that is used to select the pairs that will be collected
clear	remove the pairs from the cluster nodes

Value

Returns an object of type pairs which is a data.table. This object can be used as a regular (non-cluster) set of pairs

Examples

```

library(parallel)
data("linkexample1", "linkexample2")
cl <- makeCluster(2)

pairs <- cluster_pair(cl, linkexample1, linkexample2)
local_pairs <- cluster_collect(pairs, clear = FALSE)

compare_pairs(pairs, c("lastname", "firstname", "address", "sex"))
model <- problink_em(~ lastname + firstname + address + sex, data = pairs)
predict(model, pairs, type = "mpost", add = TRUE, binary = TRUE)
# Select pairs with a mpost > 0.5
select_threshold(pairs, "selected", "mpost", 0.5)
# Collect the selected pairs
local_pairs <- cluster_collect(pairs, "selected")

stopCluster(cl)

```

`cluster_modify_pairs` *Call a function on each of the worker nodes to modify the pairs on the node*

Description

Call a function on each of the worker nodes to modify the pairs on the node

Usage

```
cluster_modify_pairs(pairs, fun, ..., new_name = NULL)
```

Arguments

<code>pairs</code>	an object or type <code>cluster_pairs</code> as created for example by <code>cluster_pair</code> .
<code>fun</code>	a function to call on each of the worker nodes. See details on the arguments of this function.
<code>...</code>	additional arguments are passed on to <code>fun</code> .
<code>new_name</code>	name of new object to assign the pairs to on the cluster nodes.

Details

The function will have to accept the following arguments as its first three arguments:

pairs the data.table with the pairs of the worker node.
x a data.table with the portion of x present on the worker node.
y a data.table with y.

The function should either return a `data.table` with the new pairs, or `NULL`. When a `data.table` is returned this values will replace the pairs when `new_name` is missing or create new pairs in the environment `new_name`. When the function returns `NULL` it is assumed that the function modified the pairs by reference (e.g. using `pairs[,new_var := new_val]`). Note that this also means that `new_name` is ignored.

Value

Will return a `cluster_pairs` object. When `new_name` is not given it will return the input pairs invisibly. Otherwise it will return a new `cluster_pairs` object.

Examples

```
# Generate some pairs
library(parallel)
data("linkexample1", "linkexample2")
cl <- makeCluster(2)
pairs <- cluster_pair(cl, linkexample1, linkexample2)
compare_pairs(pairs, c("lastname", "firstname", "address", "sex"))

# Create a new set of pairs containing a random sample of the original
# pairs.
sample <- cluster_call(pairs, new_name = "sample", function(pairs, ...) {
  sel <- sample(nrow(pairs), round(nrow(pairs)*0.1))
  pairs[sel, ]
})

# Cleanup
stopCluster(cl)
```

cluster_pair

Generate all possible pairs using multiple processes

Description

Generates all combinations of records from `x` and `y`.

Usage

```
cluster_pair(cluster, x, y, deduplication = FALSE, name = "default")
```

Arguments

<code>cluster</code>	a cluster object as created by <code>makeCluster</code> from <code>parallel</code> or <code>makeCluster</code> from <code>snow</code> .
<code>x</code>	first <code>data.frame</code>
<code>y</code>	second <code>data.frame</code> . Ignored when <code>deduplication = TRUE</code> .
<code>deduplication</code>	generate pairs from only <code>x</code> . Ignore <code>y</code> . This is usefull for deduplication of <code>x</code> .
<code>name</code>	the name of the resulting object to create locally on the different R processes.

Details

Generating (all) pairs of the records of two data sets, is usually the first step when linking the two data sets.

`x` is split into `length{cluster}` parts which are distributed over the worker nodes. `y` is copied to each of the nodes. On the nodes then `pair` is called. The pairs are stored in the global object `reclin_env` on the nodes in the variable name. The pairs can then be further processed using functions such as `compare_pairs`, and `tabulate_patterns`. The function `cluster_collect` collects the pairs from each of the nodes.

Value

A object of type `cluster_pairs` which is a list containing the cluster and the name of the pairs object on the cluster nodes. For the pairs objects created on the nodes see the documentation of `pair`.

See Also

`cluster_pair_blocking` and `cluster_pair_minsim` are other methods to generate pairs.

Examples

```
library(parallel)
data("linkexample1", "linkexample2")
cl <- makeCluster(2)
pairs <- cluster_pair(cl, linkexample1, linkexample2)
stopCluster(cl)
```

`cluster_pair_blocking` *Generate pairs using simple blocking using multiple processes*

Description

Generates all combinations of records from `x` and `y` where the blocking variables are equal.

Usage

```
cluster_pair_blocking(
  cluster,
  x,
  y,
  on,
  deduplication = FALSE,
  name = "default"
)
```

Arguments

cluster	a cluster object as created by <code>makeCluster</code> from <code>parallel</code> or <code>makeCluster</code> from <code>snow</code> .
x	first data.frame
y	second data.frame. Ignored when <code>deduplication = TRUE</code> .
on	the variables defining the blocks or strata for which all pairs of x and y will be generated.
deduplication	generate pairs from only x. Ignore y. This is useful for deduplication of x.
name	the name of the resulting object to create locally on the different R processes.

Details

Generating (all) pairs of the records of two data sets, is usually the first step when linking the two data sets. However, this often results in a too large number of records. Therefore, blocking is usually applied.

x is split into `length{cluster}` parts which are distributed over the worker nodes. y is copied to each of the nodes. On the nodes then `pair_blocking` is called. The pairs are stored in the global object `reclin_env` on the nodes in the variable name. The pairs can then be further processed using functions such as `compare_pairs`, and `tabulate_patterns`. The function `cluster_collect` collects the pairs from each of the nodes.

Value

A object of type `cluster_pairs` which is a list containing the cluster and the name of the pairs object on the cluster nodes. For the pairs objects created on the nodes see the documentation of `pair`.

See Also

`cluster_pair` and `cluster_pair_minsim` are other methods to generate pairs.

Examples

```
library(parallel)
data("linkexample1", "linkexample2")
cl <- makeCluster(2)
pairs <- cluster_pair_blocking(cl, linkexample1, linkexample2, "postcode")
stopCluster(cl)
```

cluster_pair_minsim *Generate pairs with a minimal similarity using multiple processes*

Description

Generates all combinations of records from x and y where the blocking variables are equal.

Usage

```
cluster_pair_minsim(  
  cluster,  
  x,  
  y,  
  on,  
  minsim = 0,  
  comparators = list(default_comparator),  
  default_comparator = identical(),  
  keep_simsum = TRUE,  
  deduplication = FALSE,  
  name = "default"  
)
```

Arguments

cluster	a cluster object as created by <code>makeCluster</code> from <code>parallel</code> or <code>makeCluster</code> from <code>snow</code> .
x	first data.frame
y	second data.frame. Ignored when <code>deduplication = TRUE</code> .
on	the variables defining the blocks or strata for which all pairs of x and y will be generated.
minsim	minimal similarity score.
comparators	named list of functions with which the variables are compared. This function should accept two vectors. Function should either return a vector or a data.table with multiple columns.
default_comparator	variables for which no comparison function is defined using comparators is compared with the function <code>default_comparator</code> .
keep_simsum	add a variable <code>minsim</code> to the result with the similarity score of the pair.
deduplication	generate pairs from only x. Ignore y. This is useful for deduplication of x.
name	the name of the resulting object to create locally on the different R processes.

Details

Generating (all) pairs of the records of two data sets, is usually the first step when linking the two data sets. However, this often results in a too large number of records. `pair_minsim` will only keep pairs with a similarity score equal or larger than `minsim`. The similarity score is calculated by summing the results of the comparators for all variables of `on`.

`x` is split into `length{cluster}` parts which are distributed over the worker nodes. `y` is copied to each of the nodes. On the nodes then `cluster_pair_minsim` is called. The pairs are stored in the global object `reclin_env` on the nodes in the variable name. The pairs can then be further processed using functions such as `compare_pairs`, and `tabulate_patterns`. The function `cluster_collect` collects the pairs from each of the nodes.

Value

A object of type `cluster_pairs` which is a list containing the cluster and the name of the pairs object on the cluster nodes. For the pairs objects created on the nodes see the documentation of `pair`.

See Also

`cluster_pair` and `cluster_pair_blocking` are other methods to generate pairs.

Examples

```
library(parallel)
data("linkexample1", "linkexample2")
cl <- makeCluster(2)
# Either address or postcode has to match to keep a pair
pairs <- cluster_pair_minsim(cl, linkexample1, linkexample2,
  on = c("postcode", "address"), minsim = 1)
stopCluster(cl)
```

```
compare_pairs.cluster_pairs
```

Compare pairs on a set of variables common in both data sets

Description

Compare pairs on a set of variables common in both data sets

Usage

```
## S3 method for class 'cluster_pairs'
compare_pairs(
  pairs,
  on,
  comparators = list(default_comparator),
```

```

    default_comparator = identical(),
    new_name = NULL,
    ...
)

compare_pairs(
  pairs,
  on,
  comparators = list(default_comparator),
  default_comparator = identical(),
  ...
)

## S3 method for class 'pairs'
compare_pairs(
  pairs,
  on,
  comparators = list(default_comparator),
  default_comparator = identical(),
  x = attr(pairs, "x"),
  y = attr(pairs, "y"),
  inplace = FALSE,
  ...
)

```

Arguments

<code>pairs</code>	<code>data.table</code> with pairs. Should contain the columns <code>.x</code> and <code>.y</code> .
<code>on</code>	character vector of variables that should be compared.
<code>comparators</code>	named list of functions with which the variables are compared. This function should accept two vectors. Function should either return a vector or a <code>data.table</code> with multiple columns.
<code>default_comparator</code>	variables for which no comparison function is defined using <code>comparators</code> is compared with the function <code>default_comparator</code> .
<code>new_name</code>	name of new object to assign the pairs to on the cluster nodes.
<code>...</code>	Ignored for now
<code>x</code>	<code>data.table</code> with one half of the pairs.
<code>y</code>	<code>data.table</code> with the other half of the pairs.
<code>inplace</code>	logical indicating whether <code>pairs</code> should be modified in place. When <code>pairs</code> is large this can be more efficient.

Details

It is assumed the variables in `on` are present in both `x` and `y`. Variables with the same names are added to `pairs`. When the comparator returns a `data.table` multiple columns are added to `pairs`. The names of these columns are `variable` pasted together with the names of the `data.table` returned by comparator (separated by `"_"`).

Value

Returns the data. table pairs with one or more columns added in case of compare_pairs.pairs.

In case of compare_pairs.cluster_pairs, compare_pair.pairs is called on each cluster node and the resulting pairs are assigned to new_name in the environment reclin_env. When new_name is not given (or equal to NULL) the original pairs on the nodes are overwritten.

compare_vars.cluster_pairs

Compare pairs on given variables

Description

Compare pairs on given variables

Usage

```
## S3 method for class 'cluster_pairs'
compare_vars(
  pairs,
  variable,
  on_x = variable,
  on_y = on_x,
  comparator = identical(),
  new_name = NULL,
  ...
)

compare_vars(
  pairs,
  variable,
  on_x = variable,
  on_y = on_x,
  comparator = identical(),
  ...
)

## S3 method for class 'pairs'
compare_vars(
  pairs,
  variable,
  on_x = variable,
  on_y = on_x,
  comparator = identical(),
  x = attr(pairs, "x"),
  y = attr(pairs, "y"),
  inplace = FALSE,
```

```
    ...
  )
```

Arguments

<code>pairs</code>	<code>data.table</code> with pairs. Should contain the columns <code>.x</code> and <code>.y</code> .
<code>variable</code>	character vector with name of resulting column name that is added to pairs.
<code>on_x</code>	character vector with the column names from <code>x</code> on which to compare.
<code>on_y</code>	character vector with the column names from <code>y</code> on which to compare.
<code>comparator</code>	function with which the variables are compared. When <code>on_x</code> and <code>on_y</code> have length 1, this function should accept two vectors. Otherwise it will receive two <code>data.tables</code> . Function should either return a vector or a <code>data.table</code> with multiple columns.
<code>new_name</code>	name of new object to assign the pairs to on the cluster nodes.
<code>...</code>	Used to pass additional arguments to methods
<code>x</code>	<code>data.table</code> with one half of the pairs.
<code>y</code>	<code>data.table</code> with the other half of the pairs.
<code>inplace</code>	logical indicating whether pairs should be modified in place. When pairs is large this can be more efficient.

Details

When `comparator` returns a `data.table` multiple columns are added to pairs. The names of these columns are `variable` pasted together with the names of the `data.table` returned by `comparator` (separated by "_").

Value

Returns the `data.table` pairs with one or more columns added.

deduplicate_equivalence

Deduplication using equivalence groups

Description

Deduplication using equivalence groups

Usage

```
deduplicate_equivalence(pairs, variable, selection, x = attr(pairs, "x"))
```

Arguments

pairs	a pairs object, such as generated by pair_blocking
variable	name of the variable to create in x that will contain the group labels.
selection	a logical variable with the same length as pairs has rows, or the name of such a variable in pairs. Pairs are only selected when select is TRUE. When missing it is assumed all pairs are selected.
x	the first data set; when missing attr(pairs, "x") is used.

Value

Returns x with a variable containing the group labels. Records with the same group label (should) correspond to the same entity.

get_inspect_pairs	<i>Get a subset of pairs to inspect</i>
-------------------	-----------------------------------------

Description

Get a subset of pairs to inspect

Usage

```
get_inspect_pairs(
  pairs,
  variable,
  threshold,
  position = NULL,
  n = 11,
  x = attr(pairs, "x"),
  y = attr(pairs, "y")
)
```

Arguments

pairs	data.table with pairs.
variable	name of variable to base the selection on; should be a variable with the similarity score of the pairs.
threshold	the threshold around which to select pairs. Used when position is not given.
position	select pairs around this position (based on order of variable), e.g. position = 1 will select the pairs with the highest similarity score.
n	number of pairs to select. Pairs are selected symmetric around the threshold.
x	data.table with one half of the pairs.
y	data.table with the other half of the pairs.

Value

Returns a list with elements `pairs` with the selected pairs; `x` records from `x` corresponding to the pairs; `y` records from `y` corresponding to the pairs; `position` position of the selected pairs; `index` index of the pairs in `pairs`.

greedy	<i>Greedy one-to-one matching of pairs</i>
--------	--------------------------------------------

Description

Greedy one-to-one matching of pairs

Usage

```
greedy(x, y, weight)
```

Arguments

<code>x</code>	id's of lhs of pairs
<code>y</code>	id's of rhs of pairs
<code>weight</code>	weight of pair

Details

Pairs with the highest weight are selected as long as neither the lhs as the rhs are already selected in a pair with a higher weight.

Value

A logical vector with the same length as `x`.

identical	<i>Comparison functions</i>
-----------	-----------------------------

Description

Comparison functions

Usage

```
identical()

jaro_winkler(threshold = 0.95)

lcs(threshold = 0.8)

jaccard(threshold = 0.8)
```

Arguments

`threshold` threshold to use for the Jaro-Winkler string distance when creating a binary result.

Details

A comparison function should accept two arguments: both vectors. When the function is called with both arguments it should compare the elements in the first vector to those in the second. When called in this way, both vectors have the same length. What the function should return depends on the methods used to score the pairs. Usually the comparison functions return a similarity score with a value of 0 indicating complete difference and a value > 0 indicating similarity (often a value of 1 will indicate perfect similarity).

Some methods, such as [problink_em](#), can handle similarity scores, but also need binary values (0/FALSE = complete dissimilarity; 1/TRUE = complete similarity). In order to allow for this the comparison function is called with one argument.

When the comparison is called with one argument, it is passed the result of a previous comparison. The function should translate that result to a binary (TRUE/FALSE or 1/0) result. The result should not contain missing values.

The `jaro_winkler`, `lcs` and `jaccard` functions use the corresponding methods from [stringdist](#) except that they are transformed from a distance to a similarity score.

Value

The functions return a comparison function (see details).

Examples

```
cmp <- identical()
x <- cmp(c("john", "mary", "susan", "jack"),
        c("johan", "mary", "susanna", NA))
# Applying the comparison function to the result of the comparison results
# in a logical result, with NA's and values of FALSE set to FALSE
cmp(x)

cmp <- jaro_winkler(0.95)
x <- cmp(c("john", "mary", "susan", "jack"),
        c("johan", "mary", "susanna", NA))
# Applying the comparison function to the result of the comparison results
# in a logical result, with NA's and values below the threshold FALSE
cmp(x)
```

link*Use the selected pairs to generate a linked data set*

Description

Use the selected pairs to generate a linked data set

Usage

```
link(  
  pairs,  
  selection = NULL,  
  all = FALSE,  
  all_x = all,  
  all_y = all,  
  x = attr(pairs, "x"),  
  y = attr(pairs, "y"),  
  suffixes = c(".x", ".y"),  
  keep_from_pairs = c(".x", ".y")  
)
```

Arguments

<code>pairs</code>	a pairs object, such as generated by pair_blocking
<code>selection</code>	a logical variable with the same length as <code>pairs</code> has rows, or the name of such a variable in <code>pairs</code> . Pairs are only selected when <code>select</code> is TRUE. When missing <code>attr(pairs, "selection")</code> is used when available.
<code>all</code>	return all records from <code>x</code> and <code>y</code> ; even those that don't match.
<code>all_x</code>	return all records from <code>x</code> .
<code>all_y</code>	return all records from <code>y</code> .
<code>x</code>	the first data set; when missing <code>attr(pairs, "x")</code> is used.
<code>y</code>	the second data set; when missing <code>attr(pairs, "y")</code> is used.
<code>suffixes</code>	a character vector of length 2 specifying the suffixes to be used for making unique the names of columns in the result.
<code>keep_from_pairs</code>	character vector with names of variables in <code>pairs</code> that should be included in the output.

Details

Uses the selected pairs to link the two data sets to each other. Renames variables that are in both data sets.

Value

Returns a `data.table` containing records from `x` and `y` and pairs. Columns that occur both in `x` and `y` gain a suffix indicating from which data set they are.

linkexample1	<i>Tiny example dataset for probabilistic linkage</i>
--------------	-------------------------------------------------------

Description

Contains fictional records of 7 persons.

Format

Two data frames with resp. 6 and 5 records and 6 columns.

Details

- `id` the id of the person; this contains no errors and can be used to validate the linkage.
- `lastname` the last name of the person; contains errors.
- `firstname` the first name of the persons; contains errors.
- `address` the address; contains errors.
- `sex` the sex; contains errors and missing values.
- `postcode` the postcode; contains no errors.

match_n_to_m	<i>Force n to m matching on a set of pairs</i>
--------------	------------------------------------------------

Description

Force n to m matching on a set of pairs

Usage

```
match_n_to_m(x, y, w, n = 1, m = 1)
```

Arguments

- | | |
|----------------|-----------------------------------------------------------------------------------------------------------------|
| <code>x</code> | a vector of identifiers for each x in each pair This vector should have a unique value for each element in x. |
| <code>y</code> | a vector of identifiers for each y in each pair This vector should have a unique value for each element in y. |
| <code>w</code> | a vector with weights for each pair. The algorithm will try to maximise the total weight of the selected pairs. |
| <code>n</code> | an integer. Each element of x can be linked to at most n elements of y. |
| <code>m</code> | an integer. Each element of y can be linked to at most m elements of x. |

Details

The algorithm will try to select pairs in such a way each element of x is matched to at most n elements of y and that each element of y is matched at most m elements of x . It tries to select elements in such a way that the total weight w of the selected elements is maximised.

Value

A logical vector with the same length as x indicating the selected records.

Examples

```
d <- data.frame(x=c(1,1,1,2,2,3,3), y=c(1,2,3,4,5,6,7), w=1:7)
# One-to-one matching:
d[match_n_to_m(d$x, d$y, d$w), ]

# N-to-one matching:
d[match_n_to_m(d$x, d$y, d$w, n=999), ]

# One-to-m matching:
d[match_n_to_m(d$x, d$y, d$w, m=999), ]

# N-to-M matching, e.g. select all pairs
d[match_n_to_m(d$x, d$y, d$w, n=999, m=999), ]
```

<code>pair</code>	<i>Generate all possible pairs</i>
-------------------	------------------------------------

Description

Generates all combinations of records from x and y .

Usage

```
pair(x, y, deduplication = FALSE, add_xy = TRUE)
```

Arguments

<code>x</code>	first data.frame
<code>y</code>	second data.frame. Ignored when <code>deduplication = TRUE</code> .
<code>deduplication</code>	generate pairs from only x . Ignore y . This is useful for deduplication of x .
<code>add_xy</code>	add x and y as attributes to the returned pairs. This makes calling some subsequent operations that need x and y (such as compare_pairs) easier.

Details

Generating (all) pairs of the records of two data sets, is usually the first step when linking the two data sets.

Value

A `data.table` with two columns, `.x` and `.y`, is returned. Columns `.x` and `.y` are row numbers from `data.frames .x` and `.y` respectively.

See Also

`pair_blocking` and `pair_minsim` are other methods to generate pairs.

Examples

```
data("linkexample1", "linkexample2")
pairs <- pair(linkexample1, linkexample2)
```

pair_blocking

Generate pairs using simple blocking

Description

Generates all combinations of records from `x` and `y` where the blocking variables are equal.

Usage

```
pair_blocking(x, y, on, deduplication = FALSE, add_xy = TRUE)
```

Arguments

<code>x</code>	first <code>data.frame</code>
<code>y</code>	second <code>data.frame</code> . Ignored when <code>deduplication = TRUE</code> .
<code>on</code>	the variables defining the blocks or strata for which all pairs of <code>x</code> and <code>y</code> will be generated.
<code>deduplication</code>	generate pairs from only <code>x</code> . Ignore <code>y</code> . This is useful for deduplication of <code>x</code> .
<code>add_xy</code>	add <code>x</code> and <code>y</code> as attributes to the returned pairs. This makes calling some subsequent operations that need <code>x</code> and <code>y</code> (such as <code>compare_pairs</code>) easier.

Details

Generating (all) pairs of the records of two data sets, is usually the first step when linking the two data sets. However, this often results in a too large number of records. Therefore, blocking is usually applied.

Value

A `data.table` with two columns, `.x` and `.y`, is returned. Columns `.x` and `.y` are row numbers from `data.frames .x` and `.y` respectively.

See Also

`pair` and `pair_minsim` are other methods to generate pairs.

Examples

```
data("linkexample1", "linkexample2")
pairs <- pair_blocking(linkexample1, linkexample2, "postcode")
```

pair_minsim

Generate pairs with a minimal similarity

Description

Generates all combinations of records from `x` and `y` where the blocking variables are equal.

Usage

```
pair_minsim(
  x,
  y,
  on,
  minsim = 0,
  comparators = list(default_comparator),
  default_comparator = identical(),
  keep_simsum = TRUE,
  deduplication = FALSE,
  add_xy = TRUE
)
```

Arguments

<code>x</code>	first <code>data.frame</code>
<code>y</code>	second <code>data.frame</code> . Ignored when <code>deduplication = TRUE</code> .
<code>on</code>	the variables defining the blocks or strata for which all pairs of <code>x</code> and <code>y</code> will be generated.
<code>minsim</code>	minimal similarity score.
<code>comparators</code>	named list of functions with which the variables are compared. This function should accept two vectors. Function should either return a vector or a <code>data.table</code> with multiple columns.

default_comparator	variables for which no comparison function is defined using comparators is compares with the function default_comparator.
keep_simsum	add a variable minsim to the result with the similarity score of the pair.
deduplication	generate pairs from only x. Ignore y. This is usefull for deduplication of x.
add_xy	add x and y as attributes to the returned pairs. This makes calling some subsequent operations that need x and y (such as compare_pairs easier.

Details

Generating (all) pairs of the records of two data sets, is usually the first step when linking the two data sets. However, this often results in a too large number of records. `pair_minsim` will only keep pairs with a similarity score equal or larger than `minsim`. The similarity score is calculated by summing the results of the comparators for all variables of on.

Value

A [data.table](#) with two columns, `.x` and `.y`, is returned. Columns `.x` and `.y` are row numbers from `data.frames .x` and `.y` respectively.

See Also

[pair](#) and [pair_blocking](#) are other methods to generate pairs.

Examples

```
data("linkexample1", "linkexample2")
pairs <- pair_minsim(linkexample1, linkexample2,
  on = c("postcode", "address"), minsim = 1)
# Either address or postcode has to match to keep a pair
```

predict.problink_em *Calculate weights and probabilities for pairs*

Description

Calculate weights and probabilities for pairs

Usage

```
## S3 method for class 'problink_em'
predict(
  object,
  pairs = newdata,
  newdata = NULL,
  type = c("weights", "mpost", "probs", "all"),
```

```

    binary = FALSE,
    add = FALSE,
    comparators,
    new_name = NULL,
    ...
)

```

Arguments

object	an object of type <code>problink_em</code> as produced by <code>problink_em</code> .
pairs	a object with pairs for which to calculate weights.
newdata	an alternative name for the <code>pairs</code> argument. Specify <code>newdata</code> or <code>pairs</code> .
type	a character vector of length one specifying what to calculate. See results for more information.
binary	convert comparison vectors to binary vectors using the comparison function in <code>comparators</code> .
add	add the predictions to the original <code>pairs</code> object.
comparators	a list of comparison functions (see <code>compare_pairs</code>). When missing <code>attr(pairs, 'comparators')</code> is used.
new_name	name of new object to assign the pairs to on the cluster nodes (only relevant when <code>pairs</code> is of type <code>cluster_pairs</code>).
...	unused.

Value

When `pairs` is of type `pairs`, returns a `data.table` with either the `.x` and `.y` columns from `pairs` (when `add = FALSE`) or all columns of `pairs`. To these columns are added:

- In case of `type = "weights"` a column `weights` with the calculated weights.
- In case of `type = "mpost"` a column `mpost` with the calculated posterior probabilities (probability that pair is a match given comparison vector).
- In case of `type = "prob"` the columns `mprob` and `uprob` with the m and u-probabilites and `mpost` and `upost` with the posterior m- and u-probabilities.
- In case of `type = "all"` all of the above.

In case of `compare_pairs.cluster_pairs`, `compare_pair.pairs` is called on each cluster node and the resulting `pairs` are assigned to `new_name` in the environment `reclin_env`. When `new_name` is not given (or equal to `NULL`) the original `pairs` on the nodes are overwritten.

problink_em

*Calculate EM-estimates of m- and u-probabilities***Description**

Calculate EM-estimates of m- and u-probabilities

Usage

```
problink_em(
  formula,
  data,
  patterns,
  mprobs0 = list(0.95),
  uprobs0 = list(0.02),
  p0 = 0.05,
  tol = 1e-05,
  mprob_max = 0.999,
  uprob_min = 1e-04
)
```

Arguments

formula	a formula object with the variables for which to calculate the m- and u-probabilities. Should be of the form $\sim \text{var1} + \text{var2}$.
data	data set with pairs on which to estimate the model. Alternatively one can use the patterns argument.
patterns	table of patterns (as output by tabulate_patterns).
mprobs0, uprobs0	initial values of the m- and u-probabilities. These should be lists with numeric values. The names of the elements in the list should correspond to the names in by_x in compare_pairs .
p0	the initial estimate of the probability that a pair is a match.
tol	when the change in the m and u-probabilities is smaller than tol the algorithm is stopped.
mprob_max	maximum values of the estimated m-probabilities. Values equal to one can lead to numerical instabilities.
uprob_min	maximum values of the estimated m-probabilities. Values equal to zero can lead to numerical instabilities.

Value

Returns an object of type `problink_em`. This is a list containing the estimated `mprobs`, `uprobs` and overall linkage probability `p`. It also contains the table of comparison `patterns`.

References

Fellegi, I. and A. Sunter (1969). "A Theory for Record Linkage", *Journal of the American Statistical Association*. 64 (328): pp. 1183-1210. doi: [10.2307/2286061](https://doi.org/10.2307/2286061).

Herzog, T.N., F.J. Scheuren and W.E. Winkler (2007). *Data Quality and Record Linkage Techniques*, Springer.

Examples

```
data("linkexample1", "linkexample2")
pairs <- pair_blocking(linkexample1, linkexample2, "postcode")
pairs <- compare_pairs(pairs, c("lastname", "firstname", "address", "sex"))
model <- problink_em(~ lastname + firstname + address + sex, data = pairs)
summary(model)
```

```
select_greedy.cluster_pairs
```

Select matching pairs enforcing one-to-one linkage

Description

Select matching pairs enforcing one-to-one linkage

Usage

```
## S3 method for class 'cluster_pairs'
select_greedy(
  pairs,
  variable,
  score,
  threshold,
  preselect = NULL,
  id_x = NULL,
  id_y = NULL,
  ...
)
```

```
## S3 method for class 'cluster_pairs'
select_n_to_m(
  pairs,
  variable,
  score,
  threshold,
  preselect = NULL,
  id_x = NULL,
  id_y = NULL,
  ...
)
```

```
)

select_greedy(
  pairs,
  variable,
  score,
  threshold,
  preselect = NULL,
  id_x = NULL,
  id_y = NULL,
  ...
)

## S3 method for class 'pairs'
select_greedy(
  pairs,
  variable,
  score,
  threshold,
  preselect = NULL,
  id_x = NULL,
  id_y = NULL,
  x = attr(pairs, "x"),
  y = attr(pairs, "y"),
  inplace = FALSE,
  ...
)

select_n_to_m(
  pairs,
  variable,
  score,
  threshold,
  preselect = NULL,
  id_x = NULL,
  id_y = NULL,
  ...
)

## S3 method for class 'pairs'
select_n_to_m(
  pairs,
  variable,
  score,
  threshold,
  preselect = NULL,
  id_x = NULL,
  id_y = NULL,
```

```

    x = attr(pairs, "x"),
    y = attr(pairs, "y"),
    inplace = FALSE,
    ...
  )

```

Arguments

<code>pairs</code>	a <code>pairs</code> object, such as generated by <code>pair_blocking</code>
<code>variable</code>	the name of the new variable to create in <code>pairs</code> . This will be a logical variable with a value of <code>TRUE</code> for the selected pairs.
<code>score</code>	name of the score/weight variable of the pairs. When not given and <code>attr(pairs, "score")</code> is defined, that is used.
<code>threshold</code>	the threshold to apply. Pairs with a score above the threshold are selected.
<code>preselect</code>	a logical variable with the same length as <code>pairs</code> has rows, or the name of such a variable in <code>pairs</code> . Pairs are only selected when <code>preselect</code> is <code>TRUE</code> . This interacts with <code>threshold</code> (pairs have to be selected with both conditions).
<code>id_x</code>	a integer vector with the same length as the number of rows in <code>pairs</code> , or the name of a column in <code>x</code> . This vector should identify unique objects in <code>x</code> . When not specified it is assumed that each element in <code>x</code> is unique.
<code>id_y</code>	a integer vector with the same length as the number of rows in <code>pairs</code> , or the name of a column in <code>y</code> . This vector should identify unique objects in <code>y</code> . When not specified it is assumed that each element in <code>y</code> is unique.
<code>...</code>	Used to pass additional arguments to methods
<code>x</code>	<code>data.table</code> with one half of the pairs.
<code>y</code>	<code>data.table</code> with the other half of the pairs.
<code>inplace</code>	logical indicating whether <code>pairs</code> should be modified in place. When <code>pairs</code> is large this can be more efficient.

Details

Both methods force one-to-one matching. `select_greedy` uses a greedy algorithm that selects the first pair with the highest weight. `select_n_to_m` tries to optimise the total weight of all of the selected pairs. In general this will result in a better selection. However, `select_n_to_m` uses much more memory and is much slower and, therefore, can only be used when the number of possible pairs is not too large.

Value

Returns the `pairs` with the variable given by `variable` added. This is a logical variable indicating which pairs are selected a matches.

Examples

```

data("linkexample1", "linkexample2")
pairs <- pair_blocking(linkexample1, linkexample2, "postcode")

```

```

pairs <- compare_pairs(pairs, c("lastname", "firstname", "address", "sex"))
model <- problink_em(~ lastname + firstname + address + sex, data = pairs)
pairs <- predict(model, pairs, type = "mpost", add = TRUE, binary = TRUE)

# Select pairs with a mpost > 0.5 and force one-to-one linkage
pairs <- select_n_to_m(pairs, "ntom", "mpost", 0.5)
pairs <- select_greedy(pairs, "greedy", "mpost", 0.5)
table(pairs$ntom, pairs$greedy)

# The same example as above using a cluster;
library(parallel)
cl <- makeCluster(2)
pairs <- cluster_pair_blocking(cl, linkexample1, linkexample2, "postcode")
compare_pairs(pairs, c("lastname", "firstname", "address", "sex"))
model <- problink_em(~ lastname + firstname + address + sex, data = pairs)
predict(model, pairs, type = "mpost", add = TRUE, binary = TRUE)
# Select pairs with a mpost > 0.5 and force one-to-one linkage
# select_n_to_m and select_greedy only work on pairs that are local;
# therefore we first collect the pairs
select_threshold(pairs, "selected", "mpost", 0.5)
local_pairs <- cluster_collect(pairs, "selected")
local_pairs <- select_n_to_m(local_pairs, "ntom", "mpost", 0.5)
local_pairs <- select_greedy(local_pairs, "greedy", "mpost", 0.5)
table(local_pairs$ntom, local_pairs$greedy)

stopCluster(cl)

```

```
select_threshold.cluster_pairs
```

Select matching pairs with a score above a threshold

Description

Select matching pairs with a score above a threshold

Usage

```

## S3 method for class 'cluster_pairs'
select_threshold(pairs, variable, score, threshold, new_name = NULL, ...)

select_threshold(pairs, variable, score, threshold, ...)

## S3 method for class 'pairs'
select_threshold(pairs, variable, score, threshold, inplace = FALSE, ...)

```

Arguments

`pairs` a pairs object, such as generated by [pair_blocking](#)

variable	the name of the new variable to create in pairs. This will be a logical variable with a value of TRUE for the selected pairs.
score	name of the score/weight variable of the pairs. When not given and attr(pairs, "score") is defined, that is used.
threshold	the threshold to apply. Pairs with a score above the threshold are selected.
new_name	name of new object to assign the pairs to on the cluster nodes.
...	ignored
inplace	logical indicating whether pairs should be modified in place. When pairs is large this can be more efficient.

Value

Returns the pairs with the variable given by `variable` added. This is a logical variable indicating which pairs are selected a matches.

Examples

```
data("linkexample1", "linkexample2")
pairs <- pair_blocking(linkexample1, linkexample2, "postcode")
pairs <- compare_pairs(pairs, c("lastname", "firstname", "address", "sex"))
model <- problink_em(~ lastname + firstname + address + sex, data = pairs)
pairs <- predict(model, pairs, type = "mpost", add = TRUE, binary = TRUE)
# Select pairs with a mpost > 0.5
select_threshold(pairs, "selected", "mpost", 0.5, inplace = TRUE)

# Example using cluster;
# In general the syntax is exactly the same except for the first call to
# to cluster_pair. Note the in general `inplace = TRUE` is implied when
# working with a cluster; therefore the assignment back to pairs can be
# omitted (also not a problem if it is not).
library(parallel)
data("linkexample1", "linkexample2")
cl <- makeCluster(2)
pairs <- cluster_pair(cl, linkexample1, linkexample2)
compare_pairs(pairs, c("lastname", "firstname", "address", "sex"))
model <- problink_em(~ lastname + firstname + address + sex, data = pairs)
predict(model, pairs, type = "mpost", add = TRUE, binary = TRUE)
# Select pairs with a mpost > 0.5
# Unlike the regular pairs: inplace = TRUE is implied here
select_threshold(pairs, "selected", "mpost", 0.5)
stopCluster(cl)
```

summary.problink_em *Summarise the results from [problink_em](#)*

Description

Summarise the results from [problink_em](#)

Usage

```
## S3 method for class 'problink_em'
summary(object, ...)
```

Arguments

```
object      the problink_em object.
...         ignored;
```

Value

Returns the original object with a data.frame with the patterns and corresponding m-, u-probabilities and weights added.

```
tabulate_patterns.cluster_pairs
      Create a table of comparison patterns
```

Description

Create a table of comparison patterns

Usage

```
## S3 method for class 'cluster_pairs'
tabulate_patterns(pairs, on, comparators, complete = TRUE, ...)

tabulate_patterns(pairs, on, comparators, complete = TRUE, ...)

## S3 method for class 'pairs'
tabulate_patterns(pairs, on, comparators, complete = TRUE, ...)
```

Arguments

```
pairs      a pairs object, such as generated by pair_blocking
on         variables from pairs defining the comparison patterns. When missing names(comparators)
           is used.
comparators a list with comparison functions for each of the columns. When missing or NULL,
           the function looks for columns in pairs with a comparator attribute.
complete   add patterns that do not occur in the dataset to the result (with  $n = 0$ ).
...        passed on to other methods.
```

Details

Since comparison vectors can contain continuous numbers (usually between 0 and 1), this could result in a very large number of possible comparison vectors. Therefore, the comparison vectors are passed on to the comparators in order to threshold them. This usually results in values 0 or 1. Missing values are usually codes as 0. However, this all depends on the comparison functions used. For more information see the documentation on the [comparison functions](#).

Value

Returns a `data.frame` with all unique comparison patterns that exist in `pairs`, with a column `n` added with the number of times each pattern occurs.

Examples

```
data("linkexample1", "linkexample2")
pairs <- pair_blocking(linkexample1, linkexample2, "postcode")
pairs <- compare_pairs(pairs, c("lastname", "firstname", "address", "sex"))
tabulate_patterns(pairs)
```

town_names

Spelling variations of a set of town names

Description

Contains spelling variations found in various files of a set of town/village names. Names were selected that contain 'rdam' or 'rdm'. The correct/official names are also given. This data set can be used as an example data set for deduplication

Format

Data frames with 584 records and two columns.

Details

- `name` the name of the town/village as found in the files
- `official_name` the official/correct name

Index

- * **datasets**
 - linkexample1, 18
 - town_names, 31
- add_from_x, 2
- add_from_y (add_from_x), 2
- cluster_call, 3
- cluster_collect, 4, 7, 8, 10
- cluster_modify_pairs, 5
- cluster_pair, 3–5, 6, 8, 10
- cluster_pair_blocking, 7, 7, 10
- cluster_pair_minsim, 7, 8, 9, 10
- compare_pairs, 7, 8, 10, 19, 20, 22–24
- compare_pairs
 - (compare_pairs.cluster_pairs), 10
- compare_pairs.cluster_pairs, 10
- compare_vars
 - (compare_vars.cluster_pairs), 12
- compare_vars.cluster_pairs, 12
- comparison functions, 31
- data.table, 3, 11, 13, 14, 20–22
- deduplicate_equivalence, 13
- get_inspect_pairs, 14
- greedy, 15
- identical, 15
- jaccard (identical), 15
- jaro_winkler (identical), 15
- lcs (identical), 15
- link, 17
- linkexample1, 18
- linkexample2 (linkexample1), 18
- makeCluster, 6, 8, 9
- match_n_to_m, 18
- pair, 7, 8, 10, 19, 21, 22
- pair_blocking, 8, 14, 17, 20, 20, 22, 27, 28, 30
- pair_minsim, 20, 21, 21
- predict.problink_em, 22
- problink_em, 16, 23, 24, 29, 30
- select_greedy
 - (select_greedy.cluster_pairs), 25
- select_greedy.cluster_pairs, 25
- select_n_to_m
 - (select_greedy.cluster_pairs), 25
- select_threshold
 - (select_threshold.cluster_pairs), 28
- select_threshold.cluster_pairs, 28
- stringdist, 16
- summary.problink_em, 29
- tabulate_patterns, 7, 8, 10, 24
- tabulate_patterns
 - (tabulate_patterns.cluster_pairs), 30
- tabulate_patterns.cluster_pairs, 30
- town_names, 31