

# Package ‘csmtools’

November 3, 2016

**Type** Package

**Title** R code in use in CSM solutions

**Version** 1.1.0

**Author** McClelland Legge <McClelland.Kemp@iriworldwide.com> [aut, cre]

**Maintainer** <McClelland.Kemp@iriworldwide.com>

**Description** A collection of helpful functions that are widely or commonly used by the organization.

**License** file LICENSE

**LazyData** TRUE

**RoxygenNote** 5.0.1

**Imports** data.table,magrittr

**Suggests** testthat

## R topics documented:

csmtools-package . . . . .	2
dt_compare . . . . .	2
dt_reduce . . . . .	3
file_size_filter . . . . .	4
filter_files . . . . .	4
floor . . . . .	5
has_hive . . . . .	6
hive_datatypes . . . . .	6
hive_read . . . . .	7
hread . . . . .	7
iri_week . . . . .	8
is_defined . . . . .	9
make_compare_names . . . . .	9
ninja_load . . . . .	10

<b>Index</b>	<b>11</b>
--------------	-----------

---

csmtools-package

*csmtools: A collection of useful code*


---

## Description

csmtools is a collection of code that can be used in different projects to help prevent team members from reinventing the wheel for some things that others have already figured out.

## Details

Bug reports, additions and enhancement requests are welcomed at: <https://github.com/McClellandLegge/csmtools>

---

dt\_compare

*Merge and compare columns of data.frames (data.tables)*


---

## Description

Merge and compare columns of data.frames (data.tables)

## Usage

```
dt_compare(x, y, compare = NULL, func = `-`, precision = 6, ...)
```

## Arguments

x	A data.frame
y	A data.frame
compare	A character string or vector of shared column names
func	A binary function to compare the columns with, should be appropriate for the datatypes of the columns
precision	The precision of the comparison, is the digits argument to the <a href="#">round</a> function
...	Any arguments to the <a href="#">merge</a> function

## Value

A data.frame

## Examples

```
x <- iris[1:50,]
y <- iris[1:60,]
x$id <- seq(nrow(x))
y$id <- seq(nrow(y))
y$Sepal.Width = y$Sepal.Width + rnorm(n = nrow(y))

# can specify any arguments to 'merge'
res <- dt_compare(x, y, compare = c("Sepal.Width", "Sepal.Length"), by = "id", all.y = TRUE)
```

---

dt_reduce	<i>Apply a row-wise Reduce</i>
-----------	--------------------------------

---

## Description

Apply a row-wise Reduce

## Usage

```
dt_reduce(DT, FUN, ...)
```

## Arguments

DT	A <a href="#">data.table</a>
FUN	Any binary function
...	Quoted column names from DT

## Details

Apply a row-wise reduce for a given function on a set of a `data.table`'s columns. The main advantage of this function is that names can be passed to the function as vectors, eliminating the need to hard code differencing, etc. based on column names. Additionally, the output is specified by the user – often we want to perform a calculation and have vector output, something usually implemented with an ugly [unlist](#).

## Value

A vector  
Class will vary

## Examples

```
library("data.table")
DT <- as.data.table(head(iris))
# basic differencing
dt_reduce(DT, `-`, "Sepal.Length", "Sepal.Width")

# paste columns together row-wise
dt_reduce(DT, paste, colnames(DT))

# calculate the mean
dt_reduce(DT, `+`, "Sepal.Length", "Sepal.Width", "Petal.Length") / nrow(DT)
```

---

file_size_filter	<i>Filter files based on their size</i>
------------------	---

---

### Description

Filter files based on their size

### Usage

```
file_size_filter(x, size = 0, units = "B", include = FALSE)
```

### Arguments

x	A character vector of filenames
size	A numeric vector, the size of the file
units	A character vector specifying the units. Options are B , KB, MB and GB. Must match the length of size if specifying more than one unit.
include	A boolean, include files of size size?

### Details

Filters out the files that are less than (or less than or equal to) the size with units specified. If all files are filtered out, then a character vector of length 0 is returned.

### Value

A character vector

### Examples

```
x <- list.files(path = Sys.getenv("TEMP"), full.names = TRUE)[1]
file_size_filter(x, size = 1, units = "KB")
```

---

filter_files	<i>Filter out file names based on criteria</i>
--------------	--

---

### Description

Filter out file names based on criteria

### Usage

```
filter_files(x, size = 0, units = "B", include = FALSE, simplify = TRUE)
```

**Arguments**

x	A character vector of filenames
size	A numeric vector, the size of the file
units	A character vector specifying the units. Options are B , KB, MB and GB. Must match the length of size if specifying more than one unit.
include	A boolean, include files of size size?
simplify	A boolean, should we create a data.table from the list of output?

**Details**

Automatically excludes any non-existent items

**Value**

A list or a [data.table](#)

**Examples**

```
# single file size and unit specification
x <- list.files(path = Sys.getenv("TEMP"), full.names = TRUE)[1:100]
filter_files(x, size = 1, units = "KB")

# multiple specifications
size <- c(0, rep(1, 4))
units <- c("B", "B", "KB", "MB", "GB")
res <- filter_files(x, size, units)
res$min_file_size <- factor(res$min_file_size, levels = paste(size, units))
table(res$min_file_size)
```

---

floor

*A function to extend the functionality of the base::floor function*

---

**Description**

A function to extend the functionality of the base::floor function

**Usage**

```
floor(x, digits = 0)
```

**Arguments**

x	A numeric
digits	The number of digits to the left of the decimal place to round to

**Details**

Usually used for purchase data when you need to floor the cents

**Value**

A numeric

**Examples**

```
x <- 99.9999
floor(x, 2)
# is equivalent to:
base::floor(100 * x)
```

---

has_hive	<i>Check if the system has hive capabilities</i>
----------	--

---

**Description**

Check if the system has hive capabilities

**Usage**

```
has_hive()
```

**Details**

Checks to see if the hive binaries are in the PATH variable

**Value**

A boolean

**Examples**

```
if(has_hive()) {
  print("yes")
} else {
  print("no")
}
```

---

hive_datatypes	<i>Extract the R-datatypes for a hive table</i>
----------------	---

---

**Description**

Extract the R-datatypes for a hive table

**Usage**

```
hive_datatypes(schema, table_name)
```

**Arguments**

schema	A character string, the name of the hive schema
table_name	A character string, the name of the hive table

### Details

For now the functions converts all number-y datatypes like integer, float, decimal to numeric and both date and string types to character.

### Value

A `data.table` with columns:

- name The column name (character)
- type The R-datatype (character)

---

hive_read	<i>Read a Hive table that is stored as a text file</i>
-----------	--

---

### Description

Read a Hive table that is stored as a text file

### Usage

```
hive_read(x, ...)
```

### Arguments

- |     |  |
|-----|--|
| x   | A character string, the directory path of the hive table to read |
| ... | Additional arguments to <code>fread</code>                       |

### Details

The function allows you to assign your own additional arguments to `fread`, but it defaults the separator to pipe ("|") and adds to the `na.string` to recognize the hive default.

### Value

A `data.table`

---

hread	<i>A function to perform a read of a hive table</i>
-------	---

---

### Description

A function to perform a read of a hive table

### Usage

```
hread(table_name, schema, schema_loc, ...)
```

**Arguments**

table_name	A character string, the name of the hive table
schema	A character string, the name of the hive schema
schema_loc	A character string, the directory path of where the schema is located on the HDFS
...	Additional arguments to <a href="#">hive_read</a>

**Details**

Will automatically read all files under the directory after finding the datatypes and column names from the hive metastore. Note that the schema can have a different physical location instead of being forced to have the schema\_loc/schema naming convention.

**Value**

A [data.table](#)

**Examples**

```
## Not run:
schema_loc <- "/mapr/mapr03r/analytic_users/msmck/csm_synd_hive_schemas/csm_syndicated/"
hread("dictionary", "csm_syndicated", schema_loc)

## End(Not run)
```

---

 iri\_week

*A function to derive the IRI week for the date specified*


---

**Description**

A function to derive the IRI week for the date specified

**Usage**

```
iri_week(x, fmt = "%Y-%m-%d", ...)
```

**Arguments**

x	A date or character string
fmt	A date format
...	Additional arguments to <code>as.Date</code> including further arguments to be passed from or to other methods, including format for <code>as.character</code> and <code>as.Date</code> methods.

**Value**

A numeric

**Examples**

```
iri_week(Sys.Date())
iri_week("Dec. 12, 2016", "%b. %d, %Y")
```



---

is\_defined

*A function to perform a read of a hive table*


---

**Description**

A function to perform a read of a hive table

**Usage**

```
is_defined(..., .all = FALSE)
```

**Arguments**

...	A list of objects to test if they are null
.all	Boolean, should we return the tests for each individual element?

**Details**

Good for testing an input(s) to a function when it might be NULL

**Examples**

```
foo <- function(x = NULL, y = NULL) {
  if (is_defined(x, y)) {
    return(paste0(x, y))
  } else if (is_defined(x)) {
    return(x)
  } else if (any(is_defined(x, y, .all = TRUE))) {
    return("one is not null")
  }
}
foo(x = 1)
foo(x = 1, y = 2)
foo(y = 2)
```

---

make\_compare\_names

*Make names for comparing two data sets*


---

**Description**

Make names for comparing two data sets

**Usage**

```
make_compare_names(compare, suffixes = c(".x", ".y"), sep = "_")
```

**Arguments**

compare	The column names to compare
suffixes	The suffixes for each set to use
sep	The separator between the names, the sep argument to the <a href="#">paste</a> function

**Value**

A character vector

**Examples**

```
make_compare_names(compare = c("dollars", "units"), suffixes = c(".hive", ".sas"))
```

---

ninja_load	<i>Load packages silently</i>
------------	-------------------------------

---

**Description**

Load packages silently

**Usage**

```
ninja_load(...)
```

**Arguments**

...                    The quoted names of the packages you wish to load with deadly silence

**Examples**

```
# load some notoriously loud packages
## Not run:
ninja_load("data.table", "bit64")

## End(Not run)
```

# Index

csmttools (csmttools-package), [2](#)  
csmttools-package, [2](#)

data.table, [3](#), [5](#), [7](#), [8](#)  
dt\_compare, [2](#)  
dt\_reduce, [3](#)

file\_size\_filter, [4](#)  
filter\_files, [4](#)  
floor, [5](#)  
fread, [7](#)

has\_hive, [6](#)  
hive\_datatypes, [6](#)  
hive\_read, [7](#), [8](#)  
hread, [7](#)

iri\_week, [8](#)  
is\_defined, [9](#)

make\_compare\_names, [9](#)  
merge, [2](#)

ninja\_load, [10](#)

paste, [9](#)

round, [2](#)

unlist, [3](#)