

# Package: stream (via r-universe)

January 24, 2025

**Title** Infrastructure for Data Stream Mining

**Version** 2.0-2.1

**Date** 2024-xx-xx

**Description** A framework for data stream modeling and associated data mining tasks such as clustering and classification. The development of this package was supported in part by NSF IIS-0948893, NSF CMMI 1728612, and NIH R21HG005912. Hahsler et al (2017) <[doi:10.18637/jss.v076.i14](https://doi.org/10.18637/jss.v076.i14)>.

**License** GPL-3

**URL** <https://github.com/mhahsler/stream>

**BugReports** <https://github.com/mhahsler/stream/issues>

**Depends** R (>= 3.5.0)

**Imports** clue, cluster, clusterGeneration, dbscan (>= 1.0-0), fpc, graphics, grDevices, magrittr, MASS, methods, mlbench, proxy (>= 0.4-7), Rcpp (>= 0.11.4), rpart, stats, utils

**Suggests** animation, DBI, dplyr, knitr, rJava, RSQLite, testthat (>= 3.0.0)

**LinkingTo** BH, Rcpp

**Config/testthat/edition** 3

**Encoding** UTF-8

**Roxygen** list(markdown = TRUE)

**RoxygenNote** 7.3.2

**Repository** <https://mhahsler.r-universe.dev>

**RemoteUrl** <https://github.com/mhahsler/stream>

**RemoteRef** HEAD

**RemoteSha** a42de98fdc5271f6dbd19f8d4170538326c8780f

## Contents

agreement . . . . .	3
animate_cluster . . . . .	4
animate_data . . . . .	6
close_stream . . . . .	7
DSAggregate . . . . .	8
DSAggregate_Sample . . . . .	9
DSAggregate_Window . . . . .	10
DSC . . . . .	12
DSCClassifier . . . . .	14
DSCClassifier_SlidingWindow . . . . .	15
DSC_BICO . . . . .	16
DSC_BIRCH . . . . .	18
DSC_DBSCAN . . . . .	19
DSC_DBSTREAM . . . . .	21
DSC_DStream . . . . .	25
DSC_EA . . . . .	29
DSC_evoStream . . . . .	30
DSC_Hierarchical . . . . .	33
DSC_Kmeans . . . . .	34
DSC_Macro . . . . .	36
DSC_Micro . . . . .	37
DSC_R . . . . .	38
DSC_Reachability . . . . .	39
DSC_Sample . . . . .	41
DSC_SlidingWindow . . . . .	42
DSC_Static . . . . .	44
DSC_TwoStage . . . . .	45
DSC_Window . . . . .	47
DSD . . . . .	48
DSD_BarsAndGaussians . . . . .	49
DSD_Benchmark . . . . .	50
DSD_Cubes . . . . .	51
DSD_Gaussians . . . . .	52
DSD_Memory . . . . .	55
DSD_MG . . . . .	57
DSD_Mixture . . . . .	59
DSD_mlbenchData . . . . .	60
DSD_mlbenchGenerator . . . . .	61
DSD_NULL . . . . .	62
DSD_ReadDB . . . . .	63
DSD_ReadStream . . . . .	65
DSD_ScaleStream . . . . .	68
DSD_Target . . . . .	68
DSD_UniformNoise . . . . .	69
DSF . . . . .	70
DSFP . . . . .	72

DSF_Convolve . . . . .	73
DSF_Downsampling . . . . .	76
DSF_dplyr . . . . .	77
DSF_ExponentialMA . . . . .	79
DSF_FeatureSelection . . . . .	80
DSF_Func . . . . .	81
DSF_Scale . . . . .	82
DSOutlier . . . . .	84
DSRegressor . . . . .	85
DSRegressor_SlidingWindow . . . . .	86
DST . . . . .	87
DST_Multi . . . . .	88
DST_SlidingWindow . . . . .	89
DST_WriteStream . . . . .	92
evaluate . . . . .	93
evaluate.DSC . . . . .	94
get_assignment . . . . .	101
get_points . . . . .	103
MGC . . . . .	104
plot.DSC . . . . .	108
plot.DSD . . . . .	110
predict . . . . .	112
prune_clusters . . . . .	113
read_saveDSC . . . . .	114
recluster . . . . .	115
reset_stream . . . . .	116
stream_pipeline . . . . .	118
update . . . . .	119
write_stream . . . . .	120
<b>Index</b>	<b>122</b>

---

agreement

*Agreement-based Measures for Clustering*


---

### Description

Calculates the agreement between two partitions, typically the known actual cluster labels and the predicted cluster labels.

### Usage

```
agreement(predicted, actual, method = "cRand", na_as_cluster = TRUE)
```

## Arguments

predicted	a vector with predicted cluster labels.
actual	the known cluster labels (ground truth).
method	the used method (see <code>clue::cl_agreement()</code> ).
na_as_cluster	logical; should NA labels (noise points) be considered its own cluster?

## Details

This convenience function is an interface to `clue::cl_agreement()`. See methods in that man page for a list of available methods. A measure typically used for clustering is the corrected Rand index (also called adjusted Rand index). Numbers close to 1 indicate a very good agreement.

## References

Hornik K (2005). A CLUE for CLUster Ensembles. *Journal of Statistical Software*, 14(12). [doi:10.18637/jss.v014.i12](https://doi.org/10.18637/jss.v014.i12)

## Examples

```
# Perfect agreement (1 and 2 are just switched)
actual <- c(2, 2, 1, 3, 2, NA)
predicted <- c(1, 1, 2, 3, 1, NA)
agreement(actual, predicted)

# No agreement
predicted <- sample(predicted)
agreement(actual, predicted)
```

---

animate_cluster	<i>Animates Plots of the Clustering Process</i>
-----------------	---

---

## Description

Generates an animation of a data stream clustering process.

## Usage

```
animate_cluster(  
  dsc,  
  dsd,  
  measure = NULL,  
  horizon = 100,  
  n = 1000,  
  type = c("auto", "micro", "macro"),  
  assign = "micro",  
  assignmentMethod = c("auto", "model", "nn"),  
  excludeNoise = FALSE,
```

```

    wait = 0.1,
    plot.args = NULL,
    ...
)

```

### Arguments

dsc	a <a href="#">DSC</a>
dsd	a <a href="#">DSD</a>
measure	the evaluation measure that should be graphed below the animation (see <a href="#">evaluate_stream()</a> .)
horizon	the number of points displayed at once/used for evaluation.
n	the number of points to be plotted
type, assign, assignmentMethod, excludeNoise	are passed on to <a href="#">evaluate_stream()</a> to calculate the evaluation measure.
wait	the time interval between each frame
plot.args	a list with plotting parameters for the clusters.
...	extra arguments are added to plot.args.

### Details

Animations are recorded using the library animation and can be replayed (which gives a smoother experience since there is no more computation done) and saved in various formats (see Examples section below).

**Note:** You need to install package **animation** and its system requirements.

### Author(s)

Michael Hahsler

### See Also

[animation::ani.replay\(\)](#) for replaying and saving animations.

Other DSC: [DSC\(\)](#), [DSC\\_Macro\(\)](#), [DSC\\_Micro\(\)](#), [DSC\\_R\(\)](#), [DSC\\_SlidingWindow\(\)](#), [DSC\\_Static\(\)](#), [DSC\\_TwoStage\(\)](#), [evaluate.DSC](#), [get\\_assignment\(\)](#), [plot.DSC\(\)](#), [predict\(\)](#), [prune\\_clusters\(\)](#), [read\\_saveDSC](#), [recluster\(\)](#)

Other plot: [animate\\_data\(\)](#), [plot.DSC\(\)](#), [plot.DSD\(\)](#)

Other evaluation: [evaluate](#), [evaluate.DSC](#)

### Examples

```

if (interactive()) {
stream <- DSD_Benchmark(1)

### animate the clustering process with evaluation
### Note: we choose to exclude noise points from the evaluation
###       measure calculation, even if the algorithm would assign

```

```

###      them to a cluster.
dbstream <- DSC_DBSTREAM(r = .04, lambda = .1, gaptime = 100, Cm = 3,
  shared_density = TRUE, alpha = .2)

animate_cluster(dbstream, stream, horizon = 100, n = 5000,
  measure = "crand", type = "macro", assign = "micro",
  plot.args = list(xlim = c(0, 1), ylim = c(0, 1), shared = TRUE))
}

```

---

animate\_data

*Animates the Plotting of a Data Streams*


---

## Description

Generates an animation of a data stream.

## Usage

```
animate_data(dsd, horizon = 100, n = 1000, wait = 0.1, plot.args = NULL, ...)
```

## Arguments

dsd	a DSD object
horizon	the number of points displayed at once/used for evaluation.
n	the number of points to be plotted
wait	the time interval between each frame
plot.args	a list with plotting parameters for the clusters.
...	extra arguments are added to plot.args.

## Details

Animations are recorded using the library animation and can be replayed (which gives a smoother experience since there is no more computation done) and saved in various formats (see Examples section below).

**Note:** You need to install package **animation** and its system requirements.

## Author(s)

Michael Hahsler

## See Also

[animation::ani.replay\(\)](#) for replaying and saving animations.

Other DSD: [DSD\(\)](#), [DSD\\_BarsAndGaussians\(\)](#), [DSD\\_Benchmark\(\)](#), [DSD\\_Cubes\(\)](#), [DSD\\_Gaussians\(\)](#), [DSD\\_MG\(\)](#), [DSD\\_Memory\(\)](#), [DSD\\_Mixture\(\)](#), [DSD\\_NULL\(\)](#), [DSD\\_ReadDB\(\)](#), [DSD\\_ReadStream\(\)](#), [DSD\\_Target\(\)](#), [DSD\\_UniformNoise\(\)](#), [DSD\\_mlbenchData\(\)](#), [DSD\\_mlbenchGenerator\(\)](#), [DSF\(\)](#), [close\\_stream\(\)](#), [get\\_points\(\)](#), [plot.DSD\(\)](#), [reset\\_stream\(\)](#)

Other plot: [animate\\_cluster\(\)](#), [plot.DSC\(\)](#), [plot.DSD\(\)](#)

## Examples

```
if (interactive()) {  
  
  stream <- DSD_Benchmark(1)  
  animate_data(stream, horizon = 100, n = 5000, xlim = c(0,1), ylim = c(0,1))  
  
  ### animations can be replayed with the animation package  
  library(animation)  
  animation::ani.options(interval = .1) ## change speed  
  ani.replay()  
  
  ### animations can also be saved as HTML, animated gifs, etc.  
  saveHTML(ani.replay())  
}
```

---

close_stream	<i>Close a Data Stream</i>
--------------	----------------------------

---

## Description

Close a data stream that needs closing (e.g., a file or a connection).

## Usage

```
close_stream(dsd, ...)
```

## Arguments

dsd	An object of class a subclass of <a href="#">DSD</a> which implements a reset function.
...	further arguments.

## Details

close\_stream() is implemented for:

- [DSD](#)
- [DSD\\_ReadCSV](#)
- [DSD\\_ReadDB](#)
- [DSD\\_ReadStream](#)
- [DSF](#)
- [DST\\_WriteStream](#)

## Author(s)

Michael Hahsler

**See Also**

Other DSD: [DSD\(\)](#), [DSD\\_BarsAndGaussians\(\)](#), [DSD\\_Benchmark\(\)](#), [DSD\\_Cubes\(\)](#), [DSD\\_Gaussians\(\)](#), [DSD\\_MG\(\)](#), [DSD\\_Memory\(\)](#), [DSD\\_Mixture\(\)](#), [DSD\\_NULL\(\)](#), [DSD\\_ReadDB\(\)](#), [DSD\\_ReadStream\(\)](#), [DSD\\_Target\(\)](#), [DSD\\_UniformNoise\(\)](#), [DSD\\_mlbenchData\(\)](#), [DSD\\_mlbenchGenerator\(\)](#), [DSF\(\)](#), [animate\\_data\(\)](#), [get\\_points\(\)](#), [plot.DSD\(\)](#), [reset\\_stream\(\)](#)

---

DSAggregate

*Data Stream Aggregator Base Classes*


---

**Description**

Abstract base classes for all DSAggregate (Data Stream Aggregator) classes to aggregate streams. DSAggregate is a [DST](#) task.

**Usage**

```
DSAggregate(...)

## S3 method for class 'DSAggregate'
update(object, dsd, n = 1, return = c("nothing", "model"), ...)

## S3 method for class 'DSAggregate'
get_points(x, ...)

## S3 method for class 'DSAggregate'
get_weights(x, ...)
```

**Arguments**

...	Further arguments.
dsd	a data stream object.
n	the number of data points used for the update.
return	a character string indicating what update returns. The default is "nothing" and "model" returns the aggregated data.
x, object	a concrete implementation of DSAggregate.

**Details**

The DSAggregate class cannot be instantiated, but it serve as a base class from which other DSAggregate subclasses inherit.

Data stream operators use `update.DSAggregate()` to process new data from the [DSD](#) stream. The result of the operator can be obtained via [get\\_points\(\)](#) and [get\\_weights\(\)](#) (if available).

**Author(s)**

Michael Hahsler



**See Also**

Other DST: [DSC\(\)](#), [DSCClassifier\(\)](#), [DSOutlier\(\)](#), [DSRegressor\(\)](#), [DST\(\)](#), [DST\\_SlidingWindow\(\)](#), [DST\\_WriteStream\(\)](#), [evaluate](#), [predict\(\)](#), [stream\\_pipeline](#), [update\(\)](#)

Other DSAggregate: [DSAggregate\\_Sample\(\)](#), [DSAggregate\\_Window\(\)](#)

**Examples**

```
DSAggregate()
```

---

DSAggregate\_Sample      *Sampling from a Data Stream (Data Stream Operator)*

---

**Description**

Extracts a sample form a data stream using Reservoir Sampling.

**Usage**

```
DSAggregate_Sample(k = 100, biased = FALSE)
```

**Arguments**

k	the number of points to be sampled from the stream.
biased	if FALSE then a regular (unbiased) reservoir sampling is used. If true then the sample is biased towards keeping more recent data points (see Details section).

**Details**

If `biased = FALSE` then the reservoir sampling algorithm by McLeod and Bellhouse (1983) is used. This sampling makes sure that each data point has the same chance to be sampled. All sampled points will have a weight of 1. Note that this might not be ideal for an evolving stream since very old data points have the same chance to be in the sample as newer points.

If `bias = TRUE` then sampling prefers newer points using the modified reservoir sampling algorithm 2.1 by Aggarwal (2006). New points are always added. They replace a random point in thre reservoir with a probability of reservoir size over k. This an exponential bias function of  $2^{-lambda}$  with  $lambda = 1/k$ .

**Value**

An object of class `DSAggregate_Sample` (subclass of `DSAggregate`).

**Author(s)**

Michael Hahsler

## References

- Vitter, J. S. (1985): Random sampling with a reservoir. *ACM Transactions on Mathematical Software*, 11(1), 37-57.
- McLeod, A.I., Bellhouse, D.R. (1983): A Convenient Algorithm for Drawing a Simple Random Sample. *Applied Statistics*, 32(2), 182-184.
- Aggarwal C. (2006) On Biased Reservoir Sampling in the Presence of Stream Evolution. *International Conference on Very Large Databases (VLDB'06)*. 607-618.

## See Also

Other DSAggregate: [DSAggregate\(\)](#), [DSAggregate\\_Window\(\)](#)

## Examples

```
set.seed(1500)

stream <- DSD_Gaussians(k = 3, noise = 0.05)

sample <- DSAggregate_Sample(k = 50)
update(sample, stream, 500)
sample

head(get_points(sample))

# apply k-means clustering to the sample (data without info columns)
km <- kmeans(get_points(sample, info = FALSE), centers = 3)
plot(get_points(sample, info = FALSE))
points(km$centers, col = "red", pch = 3, cex = 2)
```

---

DSAggregate\_Window      *Sliding Window (Data Stream Operator)*

---

## Description

Implements a sliding window data stream operator which keeps a fixed amount (window length) of the most recent data points of the stream.

## Usage

```
DSAggregate_Window(horizon = 100, lambda = 0)
```

## Arguments

horizon            the window length.

lambda            decay factor damped window model. lambda = 0 means no dampening.

**Details**

If  $\lambda$  is greater than 0 then the weight uses a damped window model (Zhu and Shasha, 2002). The weight for points in the window follows  $2^{(-\lambda * t)}$  where  $t$  is the age of the point.

**Value**

An object of class DSAggregate\_Window (subclass of [DSAggregate](#)).

**Author(s)**

Michael Hahsler

**References**

Zhu, Y. and Shasha, D. (2002). StatStream: Statistical Monitoring of Thousands of Data Streams in Real Time, Intl. Conference of Very Large Data Bases (VLDB'02).

**See Also**

Other DSAggregate: [DSAggregate\(\)](#), [DSAggregate\\_Sample\(\)](#)

**Examples**

```
set.seed(1500)

## Example 1: Basic use
stream <- DSD_Gaussians(k = 3, d = 2, noise = 0.05)

window <- DSAggregate_Window(horizon = 10)
window

# update with only two points. The window is mostly empty (NA)
update(window, stream, 2)
get_points(window)

# get weights and window as a single data.frame
get_model(window)

# update window
update(window, stream, 100)
get_points(window)

## Example 2: Implement a classifier over a sliding window
window <- DSAggregate_Window(horizon = 100)

update(window, stream, 1000)

# train the classifier on the window
library(rpart)
tree <- rpart(`.class` ~ ., data = get_points(window))
tree
```

```
# predict the class for new points from the stream
new_points <- get_points(stream, n = 100, info = FALSE)
pred <- predict(tree, new_points)
plot(new_points, col = pred)
```

---

DSC

*Data Stream Clustering Base Class*


---

## Description

Abstract base classes for Data Stream Clustering (DSC). Concrete implementations are functions starting with DSC\_ (RStudio use auto-completion with Tab to select one).

## Usage

```
DSC(...)
```

```
get_centers(x, type = c("auto", "micro", "macro"), ...)
```

```
get_weights(x, type = c("auto", "micro", "macro"), scale = NULL, ...)
```

```
get_copy(x)
```

```
nclusters(x, type = c("auto", "micro", "macro"), ...)
```

```
get_microclusters(x, ...)
```

```
get_microweights(x, ...)
```

```
get_macroclusters(x, ...)
```

```
get_macroweights(x, ...)
```

## Arguments

...	further parameter
x	a DSC object.
type	Return weights of micro- or macro-clusters in x. Auto uses the class of x to decide.
scale	a range (from, to) to scale the weights. Returns by default the raw weights.

## Details

The DSC class cannot be instantiated (calling DSC() produces only a message listing the available implementations), but they serve as a base class from which other DSC classes inherit.

Data stream clustering has typically an

- **online clustering component** (see [DSC\\_Micro](#)), and an
- **offline reclustering component** (see [DSC\\_Macro](#)).

Class DSC provides several generic functions that can operate on all DSC subclasses. See Usage and Functions sections for methods. Additional, separately documented methods are:

- [update\(\)](#) adds new data points from a stream to a clustering.
- [predict\(\)](#) predicts the cluster assignment for new data points.
- [plot\(\)](#) plots cluster centers (see [plot.DSC\(\)](#)).

[get\\_centers\(\)](#) and [get\\_weights\(\)](#) are typically overwritten by subclasses of DSC.

Since DSC objects often contain external pointers, regular saving and reading operations will fail. Use [saveDSC\(\)](#) and [readDSC\(\)](#) which will serialize the objects first appropriately.

## Functions

- [get\\_centers\(\)](#): Gets the cluster centers (micro- or macro-clusters) from a DSC object.
- [get\\_weights\(\)](#): Get the weights of the clusters in the DSC (returns 1s if not implemented by the clusterer)
- [get\\_copy\(\)](#): Create a Deep Copy of a DSC Object that contain reference classes (e.g., Java data structures for MOA).
- [nclusters\(\)](#): Returns the number of micro-clusters from the DSC object.
- [get\\_microclusters\(\)](#): Used as internal interface.
- [get\\_microweights\(\)](#): Used as internal interface.
- [get\\_macroclusters\(\)](#): Used as internal interface.
- [get\\_macroweights\(\)](#): Used as internal interface.

## Author(s)

Michael Hahsler

## See Also

Other DST: [DSAggregate\(\)](#), [DSCClassifier\(\)](#), [DSOutlier\(\)](#), [DSRegressor\(\)](#), [DST\(\)](#), [DST\\_SlidingWindow\(\)](#), [DST\\_WriteStream\(\)](#), [evaluate](#), [predict\(\)](#), [stream\\_pipeline](#), [update\(\)](#)

Other DSC: [DSC\\_Macro\(\)](#), [DSC\\_Micro\(\)](#), [DSC\\_R\(\)](#), [DSC\\_SlidingWindow\(\)](#), [DSC\\_Static\(\)](#), [DSC\\_TwoStage\(\)](#), [animate\\_cluster\(\)](#), [evaluate.DSC](#), [get\\_assignment\(\)](#), [plot.DSC\(\)](#), [predict\(\)](#), [prune\\_clusters\(\)](#), [read\\_saveDSC](#), [recluster\(\)](#)

## Examples

DSC()

```
set.seed(1000)
stream <- DSD_Gaussians(k = 3, d = 2, noise = 0.05)
dstream <- DSC_DStream(gridsize = .1, gaptime = 100)
update(dstream, stream, 500)
```

```
dstream

# get micro-cluster centers
get_centers(dstream)

# get the micro-cluster weights
get_weights(dstream)

# get the number of clusters
nclusters(dstream)

# get the whole model as a data.frame
get_model(dstream)

# D-Stream also has macro-clusters
get_weights(dstream, type = "macro")
get_centers(dstream, type = "macro")

# plot the clustering result
plot(dstream, stream)
plot(dstream, stream, type = "both")

# predict macro clusters for new points (see predict())
points <- get_points(stream, n = 5)
points

predict(dstream, points, type = "macro")
```

---

DSClassifier

*Abstract Class for Data Stream Classifiers*

---

## Description

Abstract class for data stream classifiers. More implementations can be found in package **stream-MOA**.

## Usage

```
DSClassifier(...)
```

## Arguments

... Further arguments.

## Author(s)

Michael Hahsler

**See Also**

Other DST: [DSAggregate\(\)](#), [DSC\(\)](#), [DSOutlier\(\)](#), [DSRegressor\(\)](#), [DST\(\)](#), [DST\\_SlidingWindow\(\)](#), [DST\\_WriteStream\(\)](#), [evaluate](#), [predict\(\)](#), [stream\\_pipeline](#), [update\(\)](#)

Other DSClassifier: [DSClassifier\\_SlidingWindow\(\)](#)

**Examples**

```
DSClassifier()
```

---

DSClassifier\_SlidingWindow

*DSClassifier\_SlidingWindow – Data Stream Classifier Using a Sliding Window*

---

**Description**

The classifier keeps a sliding window for the stream and rebuilds a classification model at regular intervals. By default it builds a decision tree using [rpart\(\)](#).

**Usage**

```
DSClassifier_SlidingWindow(formula, model = rpart::rpart, window, rebuild, ...)
```

**Arguments**

formula	a formula for the classification problem.
model	classifier model (that has a formula interface).
window	size of the sliding window.
rebuild	interval (number of points) for rebuilding the classifier. Set rebuild to Inf to prevent automatic rebuilding. Rebuilding can be initiated manually when calling <a href="#">update()</a> .
...	additional parameters are passed on to the classifier (default is <a href="#">rpart()</a> ).

**Details**

This constructor creates classifier based on [DST\\_SlidingWindow](#). The classifier has a [update\(\)](#) and [predict\(\)](#) method.

**Value**

An object of class [DST\\_SlidingWindow](#).

**Author(s)**

Michael Hahsler

**See Also**

Other DSCClassifier: [DSCClassifier\(\)](#)

**Examples**

```
library(stream)

# create a data stream for the iris dataset
data <- iris[sample(nrow(iris)), ]
stream <- DSD_Memory(data)

# define the stream classifier.
cl <- DSCClassifier_SlidingWindow(
  Species ~ Sepal.Length + Sepal.Width + Petal.Length,
  window = 50,
  rebuild = 10
)
cl

# update the classifier with 100 points from the stream
update(cl, stream, 100)

# predict the class for the next 50 points
newdata <- get_points(stream, n = 50)
pr <- predict(cl, newdata, type = "class")
pr

table(pr, newdata$Species)

# get the tree model
get_model(cl)
```

---

DSC\_BICO

*BICO - Fast computation of k-means coresets in a data stream*

---

**Description**

Micro Clusterer. BICO maintains a tree which is inspired by the clustering tree of BIRCH. Each node in the tree represents a subset of these points. Instead of storing all points as individual objects, only the number of points, the sum and the squared sum of the subset's points are stored as key features of each subset. Points are inserted into exactly one node.

**Usage**

```
DSC_BICO(formula = NULL, k = 5, space = 10, p = 10, iterations = 10)
```



### Arguments

<code>formula</code>	NULL to use all features in the stream or a model <code>formula</code> of the form $\sim X1 + X2$ to specify the features used for clustering. Only <code>.</code> , <code>+</code> and <code>-</code> are currently supported in the formula.
<code>k</code>	number of centers.
<code>space</code>	coreset size.
<code>p</code>	number of random projections used for nearest neighbor search in first level.
<code>iterations</code>	number of repetitions for the kmeans++ procedure in the offline component.

### Details

In this implementation, the nearest neighbor search on the first level of the tree is sped up by projecting all points to random 1-d subspaces. The first estimation of the optimal clustering cost is computed in a buffer phase at the beginning of the algorithm. This implementation interfaces the original C++ implementation available here: <http://ls2-www.cs.tu-dortmund.de/grav/de/bico>. For micro-clustering, the algorithm computes the coreset of the stream. Reclustering is performed by using the kmeans++ algorithm on the coreset.

### Author(s)

R-Interface: Matthias Carnein (<Matthias.Carnein@uni-muenster.de>), Dennis Assenmacher.  
C-Implementation: Hendrik Fichtenberger, Marc Gille, Melanie Schmidt, Chris Schwiegelshohn, Christian Sohler.

### References

Hendrik Fichtenberger, Marc Gille, Melanie Schmidt, Chris Schwiegelshohn, Christian Sohler:  
BICO: BIRCH Meets Coresets for k-Means Clustering. *ESA 2013*: 481-492.

### See Also

Other DSC\_Micro: [DSC\\_BIRCH\(\)](#), [DSC\\_DBSTREAM\(\)](#), [DSC\\_DStream\(\)](#), [DSC\\_Micro\(\)](#), [DSC\\_Sample\(\)](#), [DSC\\_Window\(\)](#), [DSC\\_evoStream\(\)](#)

### Examples

```
stream <- DSD_Gaussians(k = 3, d = 2)

BICO <- DSC_BICO(k = 3, p = 10, space = 100, iterations = 10)
update(BICO, stream, n = 500)

plot(BICO, stream)
```

---

DSC\_BIRCH

*Balanced Iterative Reducing Clustering using Hierarchies*


---

### Description

Micro Clusterer. BIRCH builds a balanced tree of Clustering Features (CFs) to summarize the stream.

### Usage

```
DSC_BIRCH(
  formula = NULL,
  threshold,
  branching,
  maxLeaf,
  maxMem = 0,
  outlierThreshold = 0.25
)
```

### Arguments

formula	NULL to use all features in the stream or a model <a href="#">formula</a> of the form $\sim X1 + X2$ to specify the features used for clustering. Only <code>.</code> , <code>+</code> and <code>-</code> are currently supported in the formula.
threshold	threshold used to check whether a new data point can be absorbed or not.
branching	branching factor (maximum amount of child nodes for a non-leaf node) of the CF-Tree.
maxLeaf	maximum number of entries within a leaf node
maxMem	memory limitation for the whole CFTree in bytes. Default is 0, indicating no memory restriction.
outlierThreshold	threshold for identifying outliers when rebuilding the CF-Tree.

### Details

A CF in the balanced tree is a tuple (n, LS, SS) which represents a cluster by storing the number of elements (n), their linear sum (LS) and their squared sum (SS). Each new observation descends the tree by following its closest CF until a leaf node is reached. It is either merged into its closest leaf-CF or inserted as a new one. All leaf-CFs form the micro-clusters. Rebuilding the tree is realized by inserting all leaf-CF nodes into a new tree structure with an increased threshold.

### Author(s)

Dennis Assenmacher (<Dennis.Assenmacher@uni-muenster.de>), Matthias Carnein (<Matthias.Carnein@uni-muenste

## References

Zhang T, Ramakrishnan R and Livny M (1996), "BIRCH: An Efficient Data Clustering Method for Very Large Databases", *In Proceedings of the 1996 ACM SIGMOD International Conference on Management of Data*. Montreal, Quebec, Canada , pp. 103-114. ACM.

Zhang T, Ramakrishnan R and Livny M (1997), "BIRCH: A new data clustering algorithm and its applications", *Data Mining and Knowledge Discovery*. Vol. 1(2), pp. 141-182.

## See Also

Other DSC\_Micro: [DSC\\_BICO\(\)](#), [DSC\\_DBSTREAM\(\)](#), [DSC\\_DStream\(\)](#), [DSC\\_Micro\(\)](#), [DSC\\_Sample\(\)](#), [DSC\\_Window\(\)](#), [DSC\\_evoStream\(\)](#)

## Examples

```
stream <- DSD_Gaussians(k = 3, d = 2)

BIRCH <- DSC_BIRCH(threshold = .1, branching = 8, maxLeaf = 20)
update(BIRCH, stream, n = 500)
BIRCH

plot(BIRCH, stream)
```

---

DSC\_DBSCAN

*DBSCAN Macro-clusterer*

---

## Description

Macro Clusterer. Implements the DBSCAN algorithm for reclustering micro-clusterings.

## Usage

```
DSC_DBSCAN(
  formula = NULL,
  eps,
  MinPts = 5,
  weighted = TRUE,
  description = NULL
)
```

## Arguments

formula	NULL to use all features in the stream or a model <a href="#">formula</a> of the form $\sim X1 + X2$ to specify the features used for clustering. Only <code>.</code> , <code>+</code> and <code>-</code> are currently supported in the formula.
eps	radius of the eps-neighborhood.
MinPts	minimum number of points required in the eps-neighborhood.
weighted	logical indicating if a weighted version of DBSCAN should be used.
description	optional character string to describe the clustering method.

## Details

DBSCAN is a weighted extended version of the implementation in **fpc** where each micro-cluster center considered a pseudo point. For weighting we use in the MinPts comparison the sum of weights of the micro-cluster instead of the number.

DBSCAN first finds core points based on the number of other points in its eps-neighborhood. Then core points are joined into clusters using reachability (overlapping eps-neighborhoods).

`update()` and `recluster()` invisibly return the assignment of the data points to clusters.

**Note** that this clustering cannot be updated iteratively and every time it is used for (re)clustering, the old clustering is deleted.

## Value

An object of class DSC\_DBSCAN (a subclass of [DSC](#), [DSC\\_R](#), [DSC\\_Macro](#)).

## Author(s)

Michael Hahsler

## References

Martin Ester, Hans-Peter Kriegel, Joerg Sander, Xiaowei Xu (1996). A density-based algorithm for discovering clusters in large spatial databases with noise. In Evangelos Simoudis, Jiawei Han, Usama M. Fayyad. *Proceedings of the Second International Conference on Knowledge Discovery and Data Mining (KDD-96)*. AAAI Press. pp. 226-231.

## See Also

Other DSC\_Macro: [DSC\\_EA\(\)](#), [DSC\\_Hierarchical\(\)](#), [DSC\\_Kmeans\(\)](#), [DSC\\_Macro\(\)](#), [DSC\\_Reachability\(\)](#), [DSC\\_SlidingWindow\(\)](#)

## Examples

```
# 3 clusters with 5% noise
stream <- DSD_Gaussians(k = 3, d = 2, noise = 0.05)

# Use a moving window for "micro-clusters and recluster with DBSCAN (macro-clusters)
cl <- DSC_TwoStage(
  micro = DSC_Window(horizon = 100),
  macro = DSC_DBSCAN(eps = .05)
)

update(cl, stream, 500)
cl

plot(cl, stream)
```

---

DSC\_DBSTREAM

*DBSTREAM Clustering Algorithm*


---

### Description

Micro Clusterer with reclustering. Implements a simple density-based stream clustering algorithm that assigns data points to micro-clusters with a given radius and implements shared-density-based reclustering.

### Usage

```
DSC_DBSTREAM(
  formula = NULL,
  r,
  lambda = 0.001,
  gaptime = 1000L,
  Cm = 3,
  metric = "Euclidean",
  noise_multiplier = 1,
  shared_density = FALSE,
  alpha = 0.1,
  k = 0,
  minweight = 0
)

get_shared_density(x, use_alpha = TRUE)

change_alpha(x, alpha)

## S3 method for class 'DSC_DBSTREAM'
plot(
  x,
  dsd = NULL,
  n = 500,
  col_points = NULL,
  dim = NULL,
  method = "pairs",
  type = c("auto", "micro", "macro", "both", "none"),
  shared_density = FALSE,
  use_alpha = TRUE,
  assignment = FALSE,
  ...
)

DSOutlier_DBSTREAM(
  formula = NULL,
  r,
```

```

lambda = 0.001,
gaptime = 1000L,
Cm = 3,
metric = "Euclidean",
outlier_multiplier = 2
)

```

### Arguments

formula	NULL to use all features in the stream or a model <a href="#">formula</a> of the form $\sim X1 + X2$ to specify the features used for clustering. Only <code>.</code> , <code>+</code> and <code>-</code> are currently supported in the formula.
r	The radius of micro-clusters.
lambda	The lambda used in the fading function.
gaptime	weak micro-clusters (and weak shared density entries) are removed every <code>gaptime</code> points.
Cm	minimum weight for a micro-cluster.
metric	metric used to calculate distances.
noise_multiplier, outlier_multiplier	multiplier for radius <code>r</code> to declare noise or outliers.
shared_density	Record shared density information. If set to TRUE then shared density is used for reclustering, otherwise reachability is used (overlapping clusters with less than $r * (1 - \alpha)$ distance are clustered together).
alpha	For shared density: The minimum proportion of shared points between to clusters to warrant combining them (a suitable value for 2D data is .3). For reachability clustering it is a distance factor.
k	The number of macro clusters to be returned if macro is true.
minweight	The proportion of the total weight a macro-cluster needs to have not to be noise (between 0 and 1).
x	A DSC_DBSTREAM object to get the shared density information from.
use_alpha	only return shared density if it exceeds alpha.
dsd	a data stream object.
n	number of plots taken from the dsd to plot.
col_points	color used for plotting.
dim	an integer vector with the dimensions to plot. If NULL then for methods "pairs" and "pc" all dimensions are used and for "scatter" the first two dimensions are plotted.
method	plot method.
type	Plot micro clusters ( <code>type="micro"</code> ), macro clusters ( <code>type="macro"</code> ), both micro and macro clusters ( <code>type="both"</code> ), outliers( <code>type="outliers"</code> ), or everything together ( <code>type="all"</code> ). <code>type="auto"</code> leaves to the class of DSC to decide.
assignment	logical; show assignment area of micro-clusters.
...	further arguments are passed on to plot or pairs in graphics.

## Details

The DBSTREAM algorithm checks for each new data point in the incoming stream, if it is below the threshold value of dissimilarity value of any existing micro-clusters, and if so, merges the point with the micro-cluster. Otherwise, a new micro-cluster is created to accommodate the new data point.

Although DSC\_DBSTREAM is a micro clustering algorithm, macro clusters and weights are available.

`update()` invisibly return the assignment of the data points to clusters. The columns are `.class` with the index of the strong micro-cluster and `.mc_id` with the permanent id of the strong micro-cluster.

`plot()` for DSC\_DBSTREAM has two extra logical parameters called `assignment` and `shared_density` which show the assignment area and the shared density graph, respectively.

`predict()` can be used to assign new points to clusters. Points are assigned to a micro-cluster if they are within its assignment area (distance is less then  $r$  times `noise_multiplier`).

DSOutlier\_DBSTREAM classifies points as outlier/noise if they that cannot be assigned to a micro-cluster representing a dense region as a outlier/noise. Parameter `outlier_multiplier` specifies how far a point has to be away from a micro-cluster as a multiplier for the radius  $r$ . A larger value means that outliers have to be farther away from dense regions and thus reduce the chance of misclassifying a regular point as an outlier.

## Value

An object of class DSC\_DBSTREAM (subclass of [DSC](#), [DSC\\_R](#), [DSC\\_Micro](#)).

## Author(s)

Michael Hahsler and Matthew Bolanos

## References

Michael Hahsler and Matthew Bolanos. Clustering data streams based on shared density between micro-clusters. *IEEE Transactions on Knowledge and Data Engineering*, 28(6):1449–1461, June 2016

## See Also

Other DSC\_Micro: [DSC\\_BICO\(\)](#), [DSC\\_BIRCH\(\)](#), [DSC\\_DStream\(\)](#), [DSC\\_Micro\(\)](#), [DSC\\_Sample\(\)](#), [DSC\\_Window\(\)](#), [DSC\\_evoStream\(\)](#)

Other DSC\_TwoStage: [DSC\\_DStream\(\)](#), [DSC\\_TwoStage\(\)](#), [DSC\\_evoStream\(\)](#)

Other DSOutlier: [DSC\\_DStream\(\)](#), [DSOutlier\(\)](#)

## Examples

```
set.seed(1000)
stream <- DSD_Gaussians(k = 3, d = 2, noise = 0.05)

# create clusterer with r = .05
```

```

dbstream <- DSC_DBSTREAM(r = .05)
update(dbstream, stream, 500)
dbstream

# check micro-clusters
nclusters(dbstream)
head(get_centers(dbstream))
plot(dbstream, stream)

# plot micro-clusters with assignment area
plot(dbstream, stream, type = "none", assignment = TRUE)

# DBSTREAM with shared density
dbstream <- DSC_DBSTREAM(r = .05, shared_density = TRUE, Cm = 5)
update(dbstream, stream, 500)
dbstream

plot(dbstream, stream)
# plot the shared density graph (several options)
plot(dbstream, stream, type = "micro", shared_density = TRUE)
plot(dbstream, stream, type = "none", shared_density = TRUE, assignment = TRUE)

# see how micro and macro-clusters relate
# each micro-cluster has an entry with the macro-cluster id
# Note: unassigned micro-clusters (noise) have an NA
microToMacro(dbstream)

# do some evaluation
evaluate_static(dbstream, stream, measure = "purity")
evaluate_static(dbstream, stream, measure = "cRand", type = "macro")

# use DBSTREAM also returns the cluster assignment
# later retrieve the cluster assignments for each point)
data("iris")
dbstream <- DSC_DBSTREAM(r = 1)
cl <- update(dbstream, iris[,-5], return = "assignment")
dbstream

head(cl)

# micro-clusters
plot(iris[,-5], col = cl$class, pch = cl$class)

# macro-clusters (2 clusters since reachability cannot separate two of the three species)
plot(iris[,-5], col = microToMacro(dbstream, cl$class))

# use DBSTREAM with a formula (cluster all variables but X2)
stream <- DSD_Gaussians(k = 3, d = 4, noise = 0.05)
dbstream <- DSC_DBSTREAM(formula = ~ . - X2, r = .2)

update(dbstream, stream, 500)
get_centers(dbstream)

```



```

# use DBSTREAM for outlier detection
stream <- DSD_Gaussians(k = 3, d = 4, noise = 0.05)
outlier_detector <- DSOutlier_DBSTREAM(r = .2)

update(outlier_detector, stream, 500)
outlier_detector

plot(outlier_detector, stream)

points <- get_points(stream, 20)
points
which(is.na(predict(outlier_detector, points)))

```

---

DSC\_DStream

*D-Stream Data Stream Clustering Algorithm*


---

### Description

Micro Clusterer with reclustering. Implements the grid-based D-Stream data stream clustering algorithm.

### Usage

```

DSC_DStream(
  formula = NULL,
  gridsize,
  lambda = 0.001,
  gaptime = 1000L,
  Cm = 3,
  Cl = 0.8,
  attraction = FALSE,
  epsilon = 0.3,
  Cm2 = Cm,
  k = NULL,
  N = 0
)

get_attraction(x, relative = FALSE, grid_type = "dense", dist = FALSE)

## S3 method for class 'DSC_DStream'
plot(
  x,
  dsd = NULL,
  n = 500,
  type = c("auto", "micro", "macro", "both"),
  grid = FALSE,
  grid_type = "used",

```

```

    assignment = FALSE,
    ...
)

DSOutlier_DStream(
  formula = NULL,
  gridsize,
  lambda = 0.001,
  gaptime = 1000L,
  Cm = 3,
  Cl = 0.8,
  outlier_multiplier = 2
)

```

### Arguments

formula	NULL to use all features in the stream or a model <a href="#">formula</a> of the form $\sim X1 + X2$ to specify the features used for clustering. Only ., + and - are currently supported in the formula.
gridsize	a single number defining the size of grid cells in each dimension. For example for 3d data, the cells would be of size <code>gridsize x gridsize x gridsize</code> .
lambda	Fading constant used function to calculate the decay factor $2^{-lambda}$ . (Note: in the paper the authors use <code>lambda</code> to denote the decay factor and not the fading constant!)
gaptime	sporadic grids are removed every <code>gaptime</code> number of points.
Cm	density threshold used to detect dense grids as a proportion of the average expected density ( $Cm > 1$ ). The average density is given by the total weight of the clustering over $N$ , the number of grid cells.
Cl	density threshold to detect sporadic grids ( $0 > Cl > Cm$ ). Transitional grids have a density between <code>Cl</code> and <code>Cm</code> .
attraction	compute and store information about the attraction between adjacent grids. If <code>TRUE</code> then attraction is used to create macro-clusters, otherwise macro-clusters are created by merging adjacent dense grids.
epsilon	overlap parameter for attraction as a proportion of <code>gridsize</code> .
Cm2	threshold on attraction to join two dense grid cells (as a proportion on the average expected attraction). In the original algorithm <code>Cm2</code> is equal to <code>Cm</code> .
k	alternative to <code>Cm2</code> (not in the original algorithm). Create <code>k</code> clusters based on attraction. In case of more than <code>k</code> unconnected components, closer groups of MCs are joined.
N	Fix the number of grid cells used for the calculation of the density thresholds with <code>Cl</code> and <code>Cm</code> . If <code>N</code> is not given (0) then the algorithm tries to determine <code>N</code> from the data. Note that this means that <code>N</code> potentially increases over time and outliers might produce an extremely large value which will lead to a sudden creation of too many dense micro-clusters. The original paper assumed that <code>N</code> is known a priori.
x	DSC_DStream object to get attraction values from.

relative	calculates relative attraction (normalized by the cluster weight).
grid_type	the attraction between what grid types should be returned?
dist	make attraction symmetric and transform into a distance.
dsd	a <a href="#">DSD</a> data stream object.
n	number of plots taken from dsd to plot.
type	Plot micro clusters (type = "micro"), macro clusters (type = "macro"), both micro and macro clusters (type = "both"), outliers (type = "outliers"), or everything together (type = "all"). type = "auto" leaves to the class of DSC to decide.
grid	logical; show the D-Stream grid instead of circles for micro-clusters.
assignment	logical; show assignment area of micro-clusters.
...	further argument are passed on.
outlier_multiplier	multiplier for assignment grid width to declare outliers.

## Details

D-Stream creates an equally spaced grid and estimates the density in each grid cell using the count of points falling in the cells. Grid cells are classified based on density into dense, transitional and sporadic cells. The density is faded after every new point by a factor of  $2^{-\lambda}$ . Every gaptime number of points sporadic grid cells are removed.

For reclustering D-Stream (2007 version) merges adjacent dense grids to form macro-clusters and then assigns adjacent transitional grids to macro-clusters. This behavior is implemented as `attraction = FALSE`.

The 2009 version of the algorithm adds the concept of attraction between grids cells. If `attraction = TRUE` is used then the algorithm produces macro-clusters based on attraction between dense adjacent grids (uses `Cm2` which in the original algorithm is equal to `Cm`).

For many functions (e.g., `get_centers()`, `plot()`), D-Stream adds a parameter `grid_type` with possible values of "dense", "transitional", "sparse", "all" and "used". This only returns the selected type of grid cells. "used" includes dense and adjacent transitional cells which are used in D-Stream for reclustering. For `plot()` D-Stream also provides extra parameters "grid" and "grid\_type" to show micro-clusters as grid cells (density represented by gray values).

`DSOutlier_DStream` classifies points that do not fall into a dense grid cell as outlier/noise. Parameter `outlier_multiplier` specifies how far the point needs to be away from a dense cell to be classified as an outlier by multiplying the grid size.

Note that `DSC_DStream` currently cannot be saved to disk using `save()` or `saveRDS()`. This functionality will be added later!

## Value

An object of class `DSC_DStream` (subclass of [DSC](#), [DSC\\_R](#), [DSC\\_Micro](#)).

## Author(s)

Michael Hahsler

## References

Yixin Chen and Li Tu. 2007. Density-based clustering for real-time stream data. In *Proceedings of the 13th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD '07)*. ACM, New York, NY, USA, 133-142.

Li Tu and Yixin Chen. 2009. Stream data clustering based on grid density and attraction. *ACM Transactions on Knowledge Discovery from Data*, 3(3), Article 12 (July 2009), 27 pages.

## See Also

Other DSC\_Micro: [DSC\\_BICO\(\)](#), [DSC\\_BIRCH\(\)](#), [DSC\\_DBSTREAM\(\)](#), [DSC\\_Micro\(\)](#), [DSC\\_Sample\(\)](#), [DSC\\_Window\(\)](#), [DSC\\_evoStream\(\)](#)

Other DSC\_TwoStage: [DSC\\_DBSTREAM\(\)](#), [DSC\\_TwoStage\(\)](#), [DSC\\_evoStream\(\)](#)

Other DSOutlier: [DSC\\_DBSTREAM\(\)](#), [DSOutlier\(\)](#)

## Examples

```
stream <- DSD_BarsAndGaussians(noise = .05)
plot(stream)

dstream1 <- DSC_DStream(gridsize = 1, Cm = 1.5)
update(dstream1, stream, 1000)
dstream1

# micro-clusters (these are "used" grid cells)
nclusters(dstream1)
head(get_centers(dstream1))

# plot (DStream provides additional grid visualization)
plot(dstream1, stream)
plot(dstream1, stream, grid = TRUE)

# look only at dense grids
nclusters(dstream1, grid_type = "dense")
plot(dstream1, stream, grid = TRUE, grid_type = "dense")

# look at transitional and sparse cells
plot(dstream1, stream, grid = TRUE, grid_type = "transitional")
plot(dstream1, stream, grid = TRUE, grid_type = "sparse")

### Macro-clusters
# standard D-Stream uses reachability
nclusters(dstream1, type = "macro")
get_centers(dstream1, type = "macro")
plot(dstream1, stream, type = "macro")
evaluate_static(dstream1, stream, measure = "crand", type = "macro")

# use attraction for reclustering
dstream2 <- DSC_DStream(gridsize = 1, attraction = TRUE, Cm = 1.5)
update(dstream2, stream, 1000)
dstream2
```

```
plot(dstream2, stream, grid = TRUE)
evaluate_static(dstream2, stream, measure = "crand", type = "macro")
```

---

DSC\_EA

*Reclustering using an Evolutionary Algorithm*

---

## Description

Macro Clusterer.

## Usage

```
DSC_EA(  
  formula = NULL,  
  k,  
  generations = 2000,  
  crossoverRate = 0.8,  
  mutationRate = 0.001,  
  populationSize = 100  
)
```

## Arguments

formula	NULL to use all features in the stream or a model <a href="#">formula</a> of the form $\sim X1 + X2$ to specify the features used for clustering. Only <code>.</code> , <code>+</code> and <code>-</code> are currently supported in the formula.
k	number of macro-clusters
generations	number of EA generations performed during reclustering
crossoverRate	cross-over rate for the evolutionary algorithm
mutationRate	mutation rate for the evolutionary algorithm
populationSize	number of solutions that the evolutionary algorithm maintains

## Details

Reclustering using an evolutionary algorithm. This approach was designed for [evoStream](#) (see [DSC\\_evoStream](#)) but can also be used for other micro-clustering algorithms.

The evolutionary algorithm uses existing clustering solutions and creates small variations of them by combining and randomly modifying them. The modified solutions can yield better partitions and thus can improve the clustering over time. The evolutionary algorithm is incremental, which allows to improve existing macro-clusters instead of recomputing them every time.

## Author(s)

Matthias Carnein <Matthias.Carnein@uni-muenster.de>

## References

Carnein M. and Trautmann H. (2018), "evoStream - Evolutionary Stream Clustering Utilizing Idle Times", Big Data Research.

## See Also

Other DSC\_Macro: [DSC\\_DBSCAN\(\)](#), [DSC\\_Hierarchical\(\)](#), [DSC\\_Kmeans\(\)](#), [DSC\\_Macro\(\)](#), [DSC\\_Reachability\(\)](#), [DSC\\_SlidingWindow\(\)](#)

## Examples

```
stream <- DSD_Gaussians(k = 3, d = 2) %>% DSD_Memory(n = 1000)

## online algorithm
dbstream <- DSC_DBSTREAM(r = 0.1)

## offline algorithm (note: we use a small number of generations
##                   to make this run faster.)
EA <- DSC_EA(k = 3, generations = 100)

## create pipeline and insert observations
two <- DSC_TwoStage(dbstream, EA)
update(two, stream, n = 1000)
two

## plot result
reset_stream(stream)
plot(two, stream)

## if we have time, evaluate additional generations. This can be
## called at any time, also between observations.
two$macro$RObj$recluster(100)

## plot improved result
reset_stream(stream)
plot(two, stream)

## alternatively: do not create twostage but apply directly
reset_stream(stream)
update(dbstream, stream, n = 1000)
recluster(EA, dbstream)
reset_stream(stream)
plot(EA, stream)
```

**Description**

Micro Clusterer with reclustering. Stream clustering algorithm based on evolutionary optimization.

**Usage**

```
DSC_evoStream(
  formula = NULL,
  r,
  lambda = 0.001,
  tgap = 100,
  k = 2,
  crossoverRate = 0.8,
  mutationRate = 0.001,
  populationSize = 100,
  initializeAfter = 2 * k,
  incrementalGenerations = 1,
  reclusterGenerations = 1000
)
```

**Arguments**

formula	NULL to use all features in the stream or a model <a href="#">formula</a> of the form $\sim X1 + X2$ to specify the features used for clustering. Only ., + and - are currently supported in the formula.
r	radius threshold for micro-cluster assignment
lambda	decay rate
tgap	time-interval between outlier detection and clean-up
k	number of macro-clusters
crossoverRate	cross-over rate for the evolutionary algorithm
mutationRate	mutation rate for the evolutionary algorithm
populationSize	number of solutions that the evolutionary algorithm maintains
initializeAfter	number of micro-cluster required for the initialization of the evolutionary algorithm.
incrementalGenerations	number of EA generations performed after each observation
reclusterGenerations	number of EA generations performed during reclustering

**Details**

The online component uses a simplified version of [DBSTREAM](#) to generate micro-clusters. The micro-clusters are then incrementally reclustered using an evolutionary algorithm. Evolutionary algorithms create slight variations by combining and randomly modifying existing solutions. By iteratively selecting better solutions, an evolutionary pressure is created which improves the clustering over time. Since the evolutionary algorithm is incremental, it is possible to apply it between

observations, e.g. in the idle time of the stream. Whenever there is idle time, we can call the `recluster()` function of the reference class to improve the macro-clusters (see example). The evolutionary algorithm can also be applied as a traditional reclustering step, or a combination of both. In addition, this implementation also allows to evaluate a fixed number of generations after each observation.

### Author(s)

Matthias Carnein <Matthias.Carnein@uni-muenster.de>

### References

Carnein M. and Trautmann H. (2018), "evoStream - Evolutionary Stream Clustering Utilizing Idle Times", Big Data Research.

### See Also

Other DSC\_Micro: [DSC\\_BICO\(\)](#), [DSC\\_BIRCH\(\)](#), [DSC\\_DBSTREAM\(\)](#), [DSC\\_DStream\(\)](#), [DSC\\_Micro\(\)](#), [DSC\\_Sample\(\)](#), [DSC\\_Window\(\)](#)

Other DSC\_TwoStage: [DSC\\_DBSTREAM\(\)](#), [DSC\\_DStream\(\)](#), [DSC\\_TwoStage\(\)](#)

### Examples

```
stream <- DSD_Gaussians(k = 3, d = 2) %>% DSD_Memory(n = 500)

## init evoStream
evoStream <- DSC_evoStream(r = 0.05, k = 3,
  incrementalGenerations = 1, reclusterGenerations = 500)

## insert observations
update(evoStream, stream, n = 500)

## micro clusters
get_centers(evoStream, type = "micro")

## micro weights
get_weights(evoStream, type = "micro")

## macro clusters
get_centers(evoStream, type = "macro")

## macro weights
get_weights(evoStream, type = "macro")

## plot result
reset_stream(stream)
plot(evoStream, stream)

## if we have time, then we can evaluate additional generations.
## This can be called at any time, also between observations.
## by default, 1 generation is evaluated after each observation and
```



```

## 1000 generations during reclustering but we set it here to 500
evoStream$RObj$recluster(500)

## plot improved result
reset_stream(stream)
plot(evoStream, stream)

## get assignment of micro to macro clusters
microToMacro(evoStream)

```

DSC\_Hierarchical

*Hierarchical Micro-Cluster Reclusterer***Description**

Macro Clusterer. Implementation of hierarchical clustering to recluster a set of micro-clusters.

**Usage**

```

DSC_Hierarchical(
  formula = NULL,
  k = NULL,
  h = NULL,
  method = "complete",
  min_weight = NULL,
  description = NULL
)

```

**Arguments**

formula	NULL to use all features in the stream or a model <a href="#">formula</a> of the form $\sim X1 + X2$ to specify the features used for clustering. Only ., + and - are currently supported in the formula.
k	The number of desired clusters.
h	Height where to cut the dendrogram.
method	the agglomeration method to be used. This should be (an unambiguous abbreviation of) one of "ward", "single", "complete", "average", "mcquitty", "median" or "centroid".
min_weight	micro-clusters with a weight less than this will be ignored for reclustering.
description	optional character string to describe the clustering method.

**Details**

Please refer to [hclust\(\)](#) for more details on the behavior of the algorithm.

[update\(\)](#) and [recluster\(\)](#) invisibly return the assignment of the data points to clusters.

**Note** that this clustering cannot be updated iteratively and every time it is used for (re)clustering, the old clustering is deleted.

**Value**

A list of class `DSC`, `DSC_R`, `DSC_Macro`, and `DSC_Hierarchical`. The list contains the following items:

<code>description</code>	The name of the algorithm in the DSC object.
<code>RObj</code>	The underlying R object.

**Author(s)**

Michael Hahsler

**See Also**

Other `DSC_Macro`: `DSC_DBSCAN()`, `DSC_EA()`, `DSC_Kmeans()`, `DSC_Macro()`, `DSC_Reachability()`, `DSC_SlidingWindow()`

**Examples**

```
stream <- DSD_Gaussians(k = 3, d = 2, noise = 0.05)

# Use a moving window for "micro-clusters and recluster with HC (macro-clusters)
cl <- DSC_TwoStage(
  micro = DSC_Window(horizon = 100),
  macro = DSC_Hierarchical(h = .1, method = "single")
)

update(cl, stream, 500)
cl

plot(cl, stream)
```

---

DSC\_Kmeans

*Kmeans Macro-clusterer*

---

**Description**

Macro Clusterer. Class implements the k-means algorithm for reclusterung a set of micro-clusters.

**Usage**

```
DSC_Kmeans(
  formula = NULL,
  k,
  weighted = TRUE,
  iter.max = 10,
  nstart = 10,
  algorithm = c("Hartigan-Wong", "Lloyd", "Forgy", "MacQueen"),
  min_weight = NULL,
  description = NULL
)
```

**Arguments**

formula	NULL to use all features in the stream or a model <a href="#">formula</a> of the form $\sim X1 + X2$ to specify the features used for clustering. Only <code>.</code> , <code>+</code> and <code>-</code> are currently supported in the formula.
k	either the number of clusters, say <code>k</code> , or a set of initial (distinct) cluster centers. If a number, a random set of (distinct) rows in <code>x</code> is chosen as the initial centers.
weighted	use a weighted k-means (algorithm is ignored).
iter.max	the maximum number of iterations allowed.
nstart	if centers is a number, how many random sets should be chosen?
algorithm	character: may be abbreviated.
min_weight	micro-clusters with a weight less than this will be ignored for reclustering.
description	optional character string to describe the clustering method.

**Details**

[update\(\)](#) and [recluster\(\)](#) invisibly return the assignment of the data points to clusters.

Please refer to function [stats::kmeans\(\)](#) for more details on the algorithm.

**Note** that this clustering cannot be updated iteratively and every time it is used for (re)clustering, the old clustering is deleted.

**Value**

An object of class `DSC_Kmeans` (subclass of [DSC](#), [DSC\\_R](#), [DSC\\_Macro](#))

**Author(s)**

Michael Hahsler

**See Also**

Other `DSC_Macro`: [DSC\\_DBSCAN\(\)](#), [DSC\\_EA\(\)](#), [DSC\\_Hierarchical\(\)](#), [DSC\\_Macro\(\)](#), [DSC\\_Reachability\(\)](#), [DSC\\_SlidingWindow\(\)](#)

**Examples**

```
# 3 clusters with 5% noise
stream <- DSD_Gaussians(k = 3, d = 2, noise = 0.05)

# Use a moving window for "micro-clusters and recluster with k-means (macro-clusters)
cl <- DSC_TwoStage(
  micro = DSC_Window(horizon = 100),
  macro = DSC_Kmeans(k = 3)
)

update(cl, stream, 500)
cl

plot(cl, stream)
```

DSC\_Macro

*Abstract Class for Macro Clusterers (Offline Component)***Description**

Abstract class for all DSC Macro Clusterers which recluster micro-clusters **offline** into final clusters called macro-clusters.

**Usage**

```
DSC_Macro(...)
```

```
microToMacro(x, micro = NULL)
```

**Arguments**

...	further arguments.
x	a DSC_Macro object that also contains information about micro-clusters.
micro	A vector with micro-cluster ids. If NULL then the assignments for all micro-clusters in x are returned.

**Details**

Data stream clustering algorithms typically consists of an **online component** that creates micro-clusters (implemented as [DSC\\_Micro](#)) and and **offline components** which is used to recluster micro-clusters into final clusters called macro-clusters. The function [recluster\(\)](#) is used extract micro-clusters from a [DSC\\_Micro](#) and create macro-clusters with a DSC\_Macro.

Available clustering methods can be found in the See Also section below.

[microToMacro\(\)](#) returns the assignment of Micro-cluster IDs to Macro-cluster IDs.

For convenience, a [DSC\\_Micro](#) and DSC\_Macro can be combined using [DSC\\_TwoStage](#).

DSC\_Macro cannot be instantiated.

**Value**

A vector of the same length as `micro` with the macro-cluster ids.

**Author(s)**

Michael Hahsler

**See Also**

Other DSC\_Macro: [DSC\\_DBSCAN\(\)](#), [DSC\\_EA\(\)](#), [DSC\\_Hierarchical\(\)](#), [DSC\\_Kmeans\(\)](#), [DSC\\_Reachability\(\)](#), [DSC\\_SlidingWindow\(\)](#)

Other DSC: [DSC\(\)](#), [DSC\\_Micro\(\)](#), [DSC\\_R\(\)](#), [DSC\\_SlidingWindow\(\)](#), [DSC\\_Static\(\)](#), [DSC\\_TwoStage\(\)](#), [animate\\_cluster\(\)](#), [evaluate.DSC](#), [get\\_assignment\(\)](#), [plot.DSC\(\)](#), [predict\(\)](#), [prune\\_clusters\(\)](#), [read\\_saveDSC](#), [recluster\(\)](#)

DSC\_Micro

*Abstract Class for Micro Clusterers (Online Component)***Description**

Abstract class for all clustering methods that can operate **online** and result in a set of micro-clusters.

**Usage**

```
DSC_Micro(...)
```

**Arguments**

```
...          further arguments.
```

**Details**

Micro-clustering algorithms are data stream mining tasks [DST](#) which implement the **online component of data stream clustering**. The clustering is performed sequentially by using [update\(\)](#) to add new points from a data stream to the clustering. The result is a set of micro-clusters that can be retrieved using [get\\_clusters\(\)](#).

Available clustering methods can be found in the See Also section below.

Many data stream clustering algorithms define both, the online and an offline component to recluster micro-clusters into larger clusters called macro-clusters. This is implemented here as class [DSC\\_TwoStage](#).

DSC\_Micro cannot be instantiated.

**Author(s)**

Michael Hahsler

**See Also**

Other DSC\_Micro: [DSC\\_BICO\(\)](#), [DSC\\_BIRCH\(\)](#), [DSC\\_DBSTREAM\(\)](#), [DSC\\_DStream\(\)](#), [DSC\\_Sample\(\)](#), [DSC\\_Window\(\)](#), [DSC\\_evoStream\(\)](#)

Other DSC: [DSC\(\)](#), [DSC\\_Macro\(\)](#), [DSC\\_R\(\)](#), [DSC\\_SlidingWindow\(\)](#), [DSC\\_Static\(\)](#), [DSC\\_TwoStage\(\)](#), [animate\\_cluster\(\)](#), [evaluate.DSC](#), [get\\_assignment\(\)](#), [plot.DSC\(\)](#), [predict\(\)](#), [prune\\_clusters\(\)](#), [read\\_saveDSC](#), [recluster\(\)](#)

**Examples**

```
stream <- DSD_BarsAndGaussians(noise = .05)

# Use a DStream to create micro-clusters
dstream <- DSC_DStream(gridsize = 1, Cm = 1.5)
update(dstream, stream, 1000)
dstream
```

```
nclusters(dstream)
plot(dstream, stream, main = "micro-clusters")
```

---

DSC\_R

*Abstract Class for Implementing R-based Clusterers*


---

## Description

Abstract class for implementing R-based clusterers.

## Usage

```
DSC_R(...)

## S3 method for class 'DSC_R'
update(
  object,
  dsd,
  n = 1L,
  verbose = FALSE,
  block = 10000L,
  return = c("nothing", "assignment", "model"),
  ...
)
```

## Arguments

...	further arguments.
object	a DSC object.
dsd	a data stream object.
n	number of data points taken from the stream.
verbose	logical; show progress?
block	process blocks of data to improve speed.
return	a character string indicating what update returns. The default is "nothing". Other possible values depend on the DST. Examples are "model" and "assignment".

## Details

[DSC\\_R](#) cannot be instantiated.

### Implementing new Classes

To implement a new clusterer you need to create an S3 class with elements description and RObj. RObj needs to be a reference class with methods:

- `cluster(newdata, ...)`
- `get_microclusters(...)`

- `get_microweights(...)`
- `get_macroclusters(...)`
- `get_macroweights(...)`
- `microToMacro(micro, ...)`

See [DSC](#) for details and parameters.

DSC\_R cannot be instantiated.

### Author(s)

Michael Hahsler

### See Also

Other DSC: [DSC\(\)](#), [DSC\\_Macro\(\)](#), [DSC\\_Micro\(\)](#), [DSC\\_SlidingWindow\(\)](#), [DSC\\_Static\(\)](#), [DSC\\_TwoStage\(\)](#), [animate\\_cluster\(\)](#), [evaluate.DSC](#), [get\\_assignment\(\)](#), [plot.DSC\(\)](#), [predict\(\)](#), [prune\\_clusters\(\)](#), [read\\_saveDSC](#), [recluster\(\)](#)

---

DSC\_Reachability

*Reachability Micro-Cluster Reclusterer*

---

### Description

Macro Clusterer. Implementation of reachability clustering (based on DBSCAN's concept of reachability) to recluster a set of micro-clusters.

### Usage

```
DSC_Reachability(
  formula = NULL,
  epsilon,
  min_weight = NULL,
  description = NULL
)
```

### Arguments

<code>formula</code>	NULL to use all features in the stream or a model <a href="#">formula</a> of the form $\sim X1 + X2$ to specify the features used for clustering. Only <code>.</code> , <code>+</code> and <code>-</code> are currently supported in the formula.
<code>epsilon</code>	radius of the epsilon-neighborhood.
<code>min_weight</code>	micro-clusters with a weight less than this will be ignored for reclustering.
<code>description</code>	optional character string to describe the clustering method.

## Details

Two micro-clusters are directly reachable if they are within each other's epsilon-neighborhood (i.e., the distance between the centers is less than epsilon). Two micro-clusters are reachable if they are connected by a chain of pairwise directly reachable micro-clusters. All mutually reachable micro-clusters are put in the same cluster.

Reachability uses internally [DSC\\_Hierarchical](#) with single link.

[update\(\)](#) and [recluster\(\)](#) invisibly return the assignment of the data points to clusters.

**Note** that this clustering cannot be updated iteratively and every time it is used for (re)clustering, the old clustering is deleted.

## Value

An object of class `DSC_Reachability`. The object contains the following items:

<code>description</code>	The name of the algorithm in the DSC object.
<code>Robj</code>	The underlying R object.

## Author(s)

Michael Hahsler

## References

Martin Ester, Hans-Peter Kriegel, Joerg Sander, Xiaowei Xu (1996). A density-based algorithm for discovering clusters in large spatial databases with noise. In Evangelos Simoudis, Jiawei Han, Usama M. Fayyad. *Proceedings of the Second International Conference on Knowledge Discovery and Data Mining (KDD-96)*. AAAI Press. pp. 226-231.

## See Also

Other DSC\_Macro: [DSC\\_DBSCAN\(\)](#), [DSC\\_EA\(\)](#), [DSC\\_Hierarchical\(\)](#), [DSC\\_Kmeans\(\)](#), [DSC\\_Macro\(\)](#), [DSC\\_SlidingWindow\(\)](#)

## Examples

```
#' # 3 clusters with 5% noise
stream <- DSD_Gaussians(k = 3, d = 2, noise = 0.05)

# Use a moving window for "micro-clusters and recluster with DBSCAN (macro-clusters)
cl <- DSC_TwoStage(
  micro = DSC_Window(horizon = 100),
  macro = DSC_Reachability(eps = .05)
)

update(cl, stream, 500)
cl

plot(cl, stream)
```



---

`DSC_Sample`*Extract a Fixed-size Sample from a Data Stream*

---

## Description

Micro Clusterer. Extracts a sample form a data stream using Reservoir Sampling ([DSAggregate\\_Sample](#)). The sample is stored as a set of micro-clusters to be compatible with other data DSC stream clustering algorithms.

## Usage

```
DSC_Sample(k = 100, biased = FALSE)
```

## Arguments

<code>k</code>	the number of points to be sampled from the stream.
<code>biased</code>	if FALSE then a regular (unbiased) reservoir sampling is used. If true then the sample is biased towards keeping more recent data points (see Details section).

## Details

If `biased = FALSE` then the reservoir sampling algorithm by McLeod and Bellhouse (1983) is used. This sampling makes sure that each data point has the same chance to be sampled. All sampled points will have a weight of 1. Note that this might not be ideal for an evolving stream since very old data points have the same chance to be in the sample as newer points.

If `bias = TRUE` then sampling prefers newer points using the modified reservoir sampling algorithm 2.1 by Aggarwal (2006). New points are always added. They replace a random point in the reservoir with a probability of reservoir size over  $k$ . This an exponential bias function of  $2^{-\lambda}$  with  $\lambda = 1 / k$ .

## Value

An object of class `DSC_Sample` (subclass of `DSC`, `DSC_R`, `DSC_Micro`).

## Author(s)

Michael Hahsler

## References

- Vitter, J. S. (1985): Random sampling with a reservoir. *ACM Transactions on Mathematical Software*, 11(1), 37-57.
- McLeod, A.I., Bellhouse, D.R. (1983): A Convenient Algorithm for Drawing a Simple Random Sample. *Applied Statistics*, 32(2), 182-184.
- Aggarwal C. (2006) On Biased Reservoir Sampling in the Presence of Stream Evolution. *International Conference on Very Large Databases (VLDB'06)*. 607-618.

**See Also**

Other DSC\_Micro: [DSC\\_BICO\(\)](#), [DSC\\_BIRCH\(\)](#), [DSC\\_DBSTREAM\(\)](#), [DSC\\_DStream\(\)](#), [DSC\\_Micro\(\)](#), [DSC\\_Window\(\)](#), [DSC\\_evoStream\(\)](#)

**Examples**

```
stream <- DSD_Gaussians(k = 3, d = 2, noise = 0.05)

sample <- DSC_Sample(k = 20)
update(sample, stream, 500)
sample

# plot micro-clusters
plot(sample, stream)

# recluster the sample with k-means
kmeans <- DSC_Kmeans(k = 3)
recluster(kmeans, sample)
plot(kmeans, stream)

# sample from an evolving stream
stream <- DSD_Benchmark(1)
sample <- DSC_Sample(k = 20)
update(sample, stream, 1000)

plot(sample, stream)
# Note: the clusters move from left to right and the sample keeps many
# outdated points

# use a biased sample to keep more recent data points
stream <- DSD_Benchmark(1)
sample <- DSC_Sample(k = 20, biased = TRUE)
update(sample, stream, 1000)
plot(sample, stream)
```

---

DSC\_SlidingWindow      *DSC\_SlidingWindow – Data Stream Clusterer Using a Sliding Window*

---

**Description**

The clusterer keeps a sliding window for the stream and rebuilds a DSC clustering model at regular intervals. By default it uses [DSC\\_Kmeans](#). Other [DSC\\_Macro](#) clusterer can be used.

**Usage**

```
DSC_SlidingWindow(formula = NULL, model = DSC_Kmeans, window, rebuild, ...)
```

**Arguments**

formula	a formula for the classification problem.
model	regression model (that has a formula interface).
window	size of the sliding window.
rebuild	interval (number of points) for rebuilding the regression. Set rebuild to Inf to prevent automatic rebuilding. Rebuilding can be initiated manually when calling update().
...	additional parameters are passed on to the clusterer (default is <a href="#">DSC_Kmeans</a> ).

**Details**

This constructor creates a clusterer based on [DST\\_SlidingWindow](#). The clusterer has a `update()` and `predict()` method.

The difference to setting up a [DSC\\_TwoStage](#) is that `DSC_SlidingWindow` rebuilds the model in regular intervals, while `DSC_TwoStage` rebuilds the model on demand.

**Value**

An object of class `DST_SlidingWindow`.

**Author(s)**

Michael Hahsler

**See Also**

Other DSC: [DSC\(\)](#), [DSC\\_Macro\(\)](#), [DSC\\_Micro\(\)](#), [DSC\\_R\(\)](#), [DSC\\_Static\(\)](#), [DSC\\_TwoStage\(\)](#), [animate\\_cluster\(\)](#), [evaluate.DSC](#), [get\\_assignment\(\)](#), [plot.DSC\(\)](#), [predict\(\)](#), [prune\\_clusters\(\)](#), [read\\_saveDSC](#), [recluster\(\)](#)

Other DSC\_Macro: [DSC\\_DBSCAN\(\)](#), [DSC\\_EA\(\)](#), [DSC\\_Hierarchical\(\)](#), [DSC\\_Kmeans\(\)](#), [DSC\\_Macro\(\)](#), [DSC\\_Reachability\(\)](#)

**Examples**

```
library(stream)

stream <- DSD_Gaussians(k = 3, d = 2, noise = 0.05)

# define the stream clusterer.
cl <- DSC_SlidingWindow(
  formula = ~ . - `class`,
  k = 3,
  window = 50,
  rebuild = 10
)
cl

# update the clusterer with 100 points from the stream
```

```

update(cl, stream, 100)

# get the cluster model
cl$model$result

plot(cl$model$result)

```

---

DSC\_Static

*Create as Static Copy of a Clustering*


---

### Description

This representation cannot perform clustering anymore, but it also does not need the supporting data structures. It only stores the cluster centers and weights.

### Usage

```

DSC_Static(
  x,
  type = c("auto", "micro", "macro"),
  k_largest = NULL,
  min_weight = NULL
)

```

### Arguments

x	The clustering (a DSD object) to copy or a list with components centers (a data frame or matrix) and weights (a vector with cluster weights).
type	which clustering to copy.
k_largest	only copy the k largest (highest weight) clusters.
min_weight	only copy clusters with a weight larger or equal to min_weight.

### Value

An object of class `DSC_Static` (sub class of `DSC`, `DSC_R`). The list also contains either `DSC_Micro` or `DSC_Macro` depending on what type of clustering was copied.

### Author(s)

Michael Hahsler

### See Also

Other DSC: `DSC()`, `DSC_Macro()`, `DSC_Micro()`, `DSC_R()`, `DSC_SlidingWindow()`, `DSC_TwoStage()`, `animate_cluster()`, `evaluate.DSC`, `get_assignment()`, `plot.DSC()`, `predict()`, `prune_clusters()`, `read_saveDSC`, `recluster()`

**Examples**

```

stream <- DSD_Gaussians(k = 3, d = 2, noise = 0.05)

dstream <- DSC_DStream(gridsize = 0.05)
update(dstream, stream, 500)
dstream
plot(dstream, stream)

# create a static copy of the clustering
static <- DSC_Static(dstream)
static
plot(static, stream)

# copy only the 5 largest clusters
static2 <- DSC_Static(dstream, k_largest = 5)
static2
plot(static2, stream)

# copy all clusters with a weight of at least .3
static3 <- DSC_Static(dstream, min_weight = .3)
static3
plot(static3, stream)

# create a manual clustering
static4 <- DSC_Static(list(
  centers = data.frame(X1 = c(1, 2), X2 = c(1, 2)),
  weights = c(1, 2)),
  type = "macro")
static4
plot(static4)

```

---

DSC\_TwoStage

*TwoStage Clustering Process*


---

**Description**

Combines an **online clustering component** ([DSC\\_Micro](#)) and an **offline reclustering component** ([DSC\\_Macro](#)) into a single process.

**Usage**

```
DSC_TwoStage(micro, macro)
```

**Arguments**

micro	Clustering algorithm used in the online stage ( <a href="#">DSC_Micro</a> )
macro	Clustering algorithm used for reclustering in the offline stage ( <a href="#">DSC_Macro</a> )

**Details**

update() runs the online micro-clustering stage and only when macro cluster centers/weights are requested using `get_centers()` or `get_weights()`, then the offline stage reclustering is automatically performed.

Available clustering methods can be found in the See Also section below.

**Value**

An object of class DSC\_TwoStage (subclass of `DSC`, `DSC_Macro`) which is a named list with elements:

- description: a description of the clustering algorithms.
- micro: The `DSD` used for creating micro clusters in the online component.
- macro: The `DSD` for offline reclustering.
- state: an environment storing state information needed for reclustering.

with the two clusterers. The names are “

**Author(s)**

Michael Hahsler

**See Also**

Other DSC\_TwoStage: `DSC_DBSTREAM()`, `DSC_DStream()`, `DSC_evoStream()`

Other DSC: `DSC()`, `DSC_Macro()`, `DSC_Micro()`, `DSC_R()`, `DSC_SlidingWindow()`, `DSC_Static()`, `animate_cluster()`, `evaluate.DSC`, `get_assignment()`, `plot.DSC()`, `predict()`, `prune_clusters()`, `read_saveDSC`, `recluster()`

**Examples**

```
stream <- DSD_Gaussians(k = 3, d = 2)

# Create a clustering process that uses a window for the online stage and
# k-means for the offline stage (reclustering)
win_km <- DSC_TwoStage(
  micro = DSC_Window(horizon = 100),
  macro = DSC_Kmeans(k = 3)
)
win_km

update(win_km, stream, 200)
win_km
win_km$micro
win_km$macro

plot(win_km, stream)
evaluate_static(win_km, stream, assign = "macro")
```

---

DSC\_Window

*A sliding window from a Data Stream*

---

### Description

Interface for DSO\_Window. Represents the points in the sliding window as micro-clusters.

### Usage

```
DSC_Window(horizon = 100, lambda = 0)
```

### Arguments

horizon            the window length.  
lambda            decay factor damped window model. lambda = 0 means no dampening.

### Details

If lambda is greater than 0 then the weight uses a damped window model (Zhu and Shasha, 2002). The weight for points in the window follows  $2^{-\lambda * t}$  where  $t$  is the age of the point.

### Value

An object of class DSC\_Window (subclass of [DSC](#), [DSC\\_R](#), [DSC\\_Micro](#)).

### Author(s)

Michael Hahsler

### References

Zhu, Y. and Shasha, D. (2002). StatStream: Statistical Monitoring of Thousands of Data Streams in Real Time, *International Conference of Very Large Data Bases (VLDB'02)*.

### See Also

Other DSC\_Micro: [DSC\\_BICO\(\)](#), [DSC\\_BIRCH\(\)](#), [DSC\\_DBSTREAM\(\)](#), [DSC\\_DStream\(\)](#), [DSC\\_Micro\(\)](#), [DSC\\_Sample\(\)](#), [DSC\\_evoStream\(\)](#)

### Examples

```
stream <- DSD_Gaussians(k = 3, d = 2, noise = 0.05)

window <- DSC_Window(horizon = 100)
window

update(window, stream, 200)
window
```

```
# plot micro-clusters
plot(window, stream)

# animation for a window using a damped window model. The weight decays
# with a half-life of 25
## Not run:
window <- DSC_Window(horizon = 25, lambda = 1 / 25)
animate_cluster(window, stream, horizon = 1, n = 100, xlim = c(0, 1), ylim = c(0, 1))

## End(Not run)
```

---

DSD

*Data Stream Data Generator Base Classes*

---

## Description

Abstract base classes for DSD (Data Stream Data Generator).

## Usage

```
DSD(...)
```

```
DSD_R(...)
```

## Arguments

... further arguments.

## Details

The DSD class cannot be instantiated, but it serves as a abstract base class from which all DSD objects inherit. Implementations can be found in the See Also section below.

DSD provides common functionality like:

- `get_points()`
- `print()`
- `plot()`
- `reset_stream()` (if available)
- `close_stream()` (if needed)

DSD\_R inherits from DSD and is the abstract parent class for DSD implemented in R. To create a new R-based implementation there are only two function that needs to be implemented for a new DSD subclass called Foo would be:

1. A creator function `DSD_Foo(...)` and
2. a method `get_points.DSD_Foo(x, n = 1L)` for that class.

For details see `vignette()`



**Author(s)**

Michael Hahsler

**See Also**

Other DSD: [DSD\\_BarsAndGaussians\(\)](#), [DSD\\_Benchmark\(\)](#), [DSD\\_Cubes\(\)](#), [DSD\\_Gaussians\(\)](#), [DSD\\_MG\(\)](#), [DSD\\_Memory\(\)](#), [DSD\\_Mixture\(\)](#), [DSD\\_NULL\(\)](#), [DSD\\_ReadDB\(\)](#), [DSD\\_ReadStream\(\)](#), [DSD\\_Target\(\)](#), [DSD\\_UniformNoise\(\)](#), [DSD\\_mlbenchData\(\)](#), [DSD\\_mlbenchGenerator\(\)](#), [DSF\(\)](#), [animate\\_data\(\)](#), [close\\_stream\(\)](#), [get\\_points\(\)](#), [plot.DSD\(\)](#), [reset\\_stream\(\)](#)

**Examples**

```
DSD()

# create data stream with three clusters in 3-dimensional space
stream <- DSD_Gaussians(k = 3, d = 3)

# get points from stream
get_points(stream, n = 5)

# plotting the data (scatter plot matrix, first and third dimension, and first
# two principal components)
plot(stream)
plot(stream, dim = c(1, 3))
plot(stream, method = "pca")
```

---

DSD\_BarsAndGaussians *Data Stream Generator for Bars and Gaussians*

---

**Description**

A data stream generator which creates the shape of two bars and two Gaussians clusters with different density.

**Usage**

```
DSD_BarsAndGaussians(angle = NULL, noise = 0)
```

**Arguments**

**angle** rotation in degrees. NULL will produce a random rotation.  
**noise** The amount of noise that should be added to the output.

**Value**

Returns a DSD\_BarsAndGaussians object.

**Author(s)**

Michael Hahsler

**See Also****DSD**

Other DSD: [DSD\(\)](#), [DSD\\_Benchmark\(\)](#), [DSD\\_Cubes\(\)](#), [DSD\\_Gaussians\(\)](#), [DSD\\_MG\(\)](#), [DSD\\_Memory\(\)](#), [DSD\\_Mixture\(\)](#), [DSD\\_NULL\(\)](#), [DSD\\_ReadDB\(\)](#), [DSD\\_ReadStream\(\)](#), [DSD\\_Target\(\)](#), [DSD\\_UniformNoise\(\)](#), [DSD\\_mlbenchData\(\)](#), [DSD\\_mlbenchGenerator\(\)](#), [DSF\(\)](#), [animate\\_data\(\)](#), [close\\_stream\(\)](#), [get\\_points\(\)](#), [plot.DSD\(\)](#), [reset\\_stream\(\)](#)

**Examples**

```
# create data stream with three clusters in 2D
stream <- DSD_BarsAndGaussians(noise = 0.1)

get_points(stream, n = 10)
plot(stream)
```

---

DSD\_Benchmark

*Data Stream Generator for Dynamic Data Stream Benchmarks*

---

**Description**

A data stream generator that generates several dynamic streams intended to be benchmarks to compare data stream clustering algorithms. The benchmarks can be used to test if a clustering algorithm can follow moving clusters, and merging and separating clusters.

**Usage**

```
DSD_Benchmark(i = 1)
```

**Arguments**

*i* integer; the number of the benchmark.

**Details**

Currently available benchmarks are:

- 1: two tight clusters moving across the data space with noise and intersect in the middle.
- 2: two clusters are located in two corners of the data space. A third cluster moves between the two clusters forth and back.

The benchmarks are created using [DSD\\_MG](#).

**Value**

Returns a [DSD](#) object.

**Author(s)**

Michael Hahsler

**See Also**

Other DSD: [DSD\(\)](#), [DSD\\_BarsAndGaussians\(\)](#), [DSD\\_Cubes\(\)](#), [DSD\\_Gaussians\(\)](#), [DSD\\_MG\(\)](#), [DSD\\_Memory\(\)](#), [DSD\\_Mixture\(\)](#), [DSD\\_NULL\(\)](#), [DSD\\_ReadDB\(\)](#), [DSD\\_ReadStream\(\)](#), [DSD\\_Target\(\)](#), [DSD\\_UniformNoise\(\)](#), [DSD\\_mlbenchData\(\)](#), [DSD\\_mlbenchGenerator\(\)](#), [DSF\(\)](#), [animate\\_data\(\)](#), [close\\_stream\(\)](#), [get\\_points\(\)](#), [plot.DSD\(\)](#), [reset\\_stream\(\)](#)

**Examples**

```
stream <- DSD_Benchmark(i = 1)
get_points(stream, n = 5)

## Not run:
stream <- DSD_Benchmark(i = 1)
animate_data(stream, n = 10000, horizon = 100, xlim = c(0, 1), ylim = c(0, 1))

stream <- DSD_Benchmark(i = 2)
animate_data(stream, n = 10000, horizon = 100, xlim = c(0, 1), ylim = c(0, 1))

## End(Not run)
```

---

DSD\_Cubes

*Static Cubes Data Stream Generator*


---

**Description**

A data stream generator that produces a data stream with static (hyper) cubes filled uniformly with data points.

**Usage**

```
DSD_Cubes(k = 2, d = 2, center, size, p, noise = 0, noise_range)
```

**Arguments**

k	Determines the number of clusters.
d	Determines the number of dimensions.
center	A matrix of means for each dimension of each cluster.
size	A k times d matrix with the cube dimensions.
p	A vector of probabilities that determines the likelihood of generated a data point from a particular cluster.
noise	Noise probability between 0 and 1. Noise is uniformly distributed within noise range (see below).
noise_range	A matrix with d rows and 2 columns. The first column contains the minimum values and the second column contains the maximum values for noise.

**Value**

Returns a DSD\_Cubes object (subclass of [DSD\\_R](#), [DSD](#)).

**Author(s)**

Michael Hahsler

**See Also**

Other DSD: [DSD\(\)](#), [DSD\\_BarsAndGaussians\(\)](#), [DSD\\_Benchmark\(\)](#), [DSD\\_Gaussians\(\)](#), [DSD\\_MG\(\)](#), [DSD\\_Memory\(\)](#), [DSD\\_Mixture\(\)](#), [DSD\\_NULL\(\)](#), [DSD\\_ReadDB\(\)](#), [DSD\\_ReadStream\(\)](#), [DSD\\_Target\(\)](#), [DSD\\_UniformNoise\(\)](#), [DSD\\_mlbenchData\(\)](#), [DSD\\_mlbenchGenerator\(\)](#), [DSF\(\)](#), [animate\\_data\(\)](#), [close\\_stream\(\)](#), [get\\_points\(\)](#), [plot.DSD\(\)](#), [reset\\_stream\(\)](#)

**Examples**

```
# create data stream with three clusters in 3D
stream <- DSD_Cubes(k = 3, d = 3, noise = 0.05)

get_points(stream, n = 5)

plot(stream)
```

---

DSD\_Gaussians

*Mixture of Gaussians Data Stream Generator*

---

**Description**

A data stream generator that produces a data stream with a mixture of static Gaussians.

**Usage**

```
DSD_Gaussians(
  k = 3,
  d = 2,
  p,
  mu,
  sigma,
  variance_limit = c(0.001, 0.002),
  separation = 6,
  space_limit = c(0, 1),
  noise = 0,
  noise_limit = space_limit,
  noise_separation = 3,
  separation_type = c("Euclidean", "Mahalanobis"),
  verbose = FALSE
)
```

**Arguments**

k	Determines the number of clusters.
d	Determines the number of dimensions.
p	A vector of probabilities that determines the likelihood of generated a data point from a particular cluster.
mu	A matrix of means for each dimension of each cluster.
sigma	A list of length k of covariance matrices.
variance_limit	Lower and upper limit for the randomly generated variance when creating cluster covariance matrices.
separation	Minimum separation distance between clusters (measured in standard deviations according to separation_type).
space_limit	Defines the space bounds. All constructs are generated inside these bounds. For clusters this means that their centroids must be within these space bounds.
noise	Noise probability between 0 and 1. Noise is uniformly distributed within noise range (see below).
noise_limit	A matrix with d rows and 2 columns. The first column contains the minimum values and the second column contains the maximum values for noise.
noise_separation	Minimum separation distance between cluster centers and noise points (measured in standard deviations according to separation_type). 0 means separation is ignored.
separation_type	The type of the separation distance calculation. It can be either Euclidean distance or Mahalanobis distance.
verbose	Report cluster and outlier generation process.

**Details**

DSD\_Gaussians creates a mixture of k static clusters in a d-dimensional space. The cluster centers mu and the covariance matrices sigma can be supplied or will be randomly generated. The probability vector p defines for each cluster the probability that the next data point will be chosen from it (defaults to equal probability). Separation between generated clusters (and outliers; see below) can be imposed by using Euclidean or Mahalanobis distance, which is controlled by the separation\_type parameter. Separation value then is supplied in the separation parameter. The generation method is similar to the one suggested by Jain and Dubes (1988).

Noise points which are uniformly chosen from noise\_limit can be added.

Outlier points can be added. The outlier spatial positions predefined\_outlier\_space\_positions and the outlier stream positions predefined\_outlier\_stream\_positions can be supplied or will be randomly generated. Cluster and outlier separation distance is determined by and outlier\_virtual\_variance parameters. The outlier virtual variance defines an empty space around outliers, which separates them from their surrounding. Unlike noise, outliers are data points of interest for end-users, and the goal of outlier detectors is to find them in data streams. For more details, read the "Introduction to **stream**" vignette.

**Value**

Returns a object of class `DSD_Gaussian` (subclass of `DSD_R`, `DSD`).

**Author(s)**

Michael Hahsler

**References**

Jain and Dubes (1988) Algorithms for clustering data, Prentice-Hall, Inc., Upper Saddle River, NJ, USA.

**See Also**

Other DSD: `DSD()`, `DSD_BarsAndGaussians()`, `DSD_Benchmark()`, `DSD_Cubes()`, `DSD_MG()`, `DSD_Memory()`, `DSD_Mixture()`, `DSD_NULL()`, `DSD_ReadDB()`, `DSD_ReadStream()`, `DSD_Target()`, `DSD_UniformNoise()`, `DSD_mlbenchData()`, `DSD_mlbenchGenerator()`, `DSF()`, `animate_data()`, `close_stream()`, `get_points()`, `plot.DSD()`, `reset_stream()`

**Examples**

```
# Example 1: create data stream with three clusters in 3-dimensional data space
#           with 5 times sqrt(variance_limit) separation.
set.seed(1)
stream1 <- DSD_Gaussians(k = 3, d = 3)
stream1

get_points(stream1, n = 5)
plot(stream1, xlim = c(0, 1), ylim = c(0, 1))
```

```
# Example 2: create data stream with specified cluster positions,
# 5% noise in a given bounding box and
# with different densities (1 to 9 between the two clusters)
stream2 <- DSD_Gaussians(k = 2, d = 2,
  mu = rbind(c(-.5, -.5), c(.5, .5)),
  p = c(.1, .9),
  variance_limit = c(0.02, 0.04),
  noise = 0.05,
  noise_limit = rbind(c(-1, 1), c(-1, 1)))
```

```
get_points(stream2, n = 5)
plot(stream2, xlim = c(-1, 1), ylim = c(-1, 1))
```

```
# Example 3: create 4 clusters and noise separated by a Mahalanobis
# distance. Distance to noise is increased to 6 standard deviations to make them
# easier detectable outliers.
stream3 <- DSD_Gaussians(k = 4, d = 2,
  separation_type = "Mahalanobis",
  space_limit = c(5, 20),
```

```

variance_limit = c(1, 2),
noise = 0.05,
noise_limit = c(0, 25),
noise_separation = 6
)
plot(stream3)

```

---

DSD\_Memory

*A Data Stream Interface for Data Stored in Memory*


---

## Description

This class provides a data stream interface for data stored in memory as matrix-like objects (including data frames). All or a portion of the stored data can be replayed several times.

## Usage

```

DSD_Memory(
  x,
  n,
  k = NA,
  outofpoints = c("warn", "ignore", "stop"),
  loop = FALSE,
  description = NULL
)

```

## Arguments

x	A matrix-like object containing the data. If x is a DSD object then a data frame for n data points from this DSD is created.
n	Number of points used if x is a DSD object. If x is a matrix-like object then n is ignored.
k	Optional: The known number of clusters in the data
outofpoints	Action taken if less than n data points are available. The default is to return the available data points with a warning. Other supported actions are: <ul style="list-style-type: none"> <li>• warn: return the available points (maybe an empty data.frame) with a warning.</li> <li>• ignore: silently return the available points.</li> <li>• stop: stop with an error.</li> </ul>
loop	Should the stream start over when it reaches the end?
description	character string with a description.

## Details

In addition to regular data.frames other matrix-like objects that provide subsetting with the bracket operator can be used. This includes `ffdf` (large data.frames stored on disk) from package **ff** and `big.matrix` from **bigmemory**.

**Reading the whole stream** By using `n = -1` in `get_points()`, the whole stream is returned.

**Value**

Returns a DSD\_Memory object (subclass of [DSD\\_R](#), [DSD](#)).

**Author(s)**

Michael Hahsler

**See Also**

Other DSD: [DSD\(\)](#), [DSD\\_BarsAndGaussians\(\)](#), [DSD\\_Benchmark\(\)](#), [DSD\\_Cubes\(\)](#), [DSD\\_Gaussians\(\)](#), [DSD\\_MG\(\)](#), [DSD\\_Mixture\(\)](#), [DSD\\_NULL\(\)](#), [DSD\\_ReadDB\(\)](#), [DSD\\_ReadStream\(\)](#), [DSD\\_Target\(\)](#), [DSD\\_UniformNoise\(\)](#), [DSD\\_mlbenchData\(\)](#), [DSD\\_mlbenchGenerator\(\)](#), [DSF\(\)](#), [animate\\_data\(\)](#), [close\\_stream\(\)](#), [get\\_points\(\)](#), [plot.DSD\(\)](#), [reset\\_stream\(\)](#)

**Examples**

```
# Example 1: store 1000 points from a stream
stream <- DSD_Gaussians(k = 3, d = 2)
replayer <- DSD_Memory(stream, k = 3, n = 1000)
replayer
plot(replayer)

# creating 2 clusterers of different algorithms
dsc1 <- DSC_DBSTREAM(r = 0.1)
dsc2 <- DSC_DStream(gridsize = 0.1, Cm = 1.5)

# clustering the same data in 2 DSC objects
reset_stream(replayer) # resetting the replayer to the first position
update(dsc1, replayer, 500)
reset_stream(replayer)
update(dsc2, replayer, 500)

# plot the resulting clusterings
reset_stream(replayer)
plot(dsc1, replayer, main = "DBSTREAM")
reset_stream(replayer)
plot(dsc2, replayer, main = "D-Stream")

# Example 2: use a data.frame to create a stream (3rd col. contains the assignment)
df <- data.frame(x = runif(100), y = runif(100),
  .class = sample(1:3, 100, replace = TRUE))

# add some outliers
out <- runif(100) > .95
df[['.outlier']] <- out
df[['.class']] <- NA
head(df)

stream <- DSD_Memory(df)
stream
```



```

reset_stream(stream)
get_points(stream, n = 5)

# get the remaining points
rest <- get_points(stream, n = -1)
nrow(rest)

# plot all available points with n = -1
reset_stream(stream)
plot(stream, n = -1)

```

---

DSD\_MG

*DSD Moving Generator*


---

### Description

Creates an evolving DSD that consists of several [MGC](#), each representing a moving cluster.

### Usage

```

DSD_MG(dimension = 2, ..., labels = NULL, description = NULL)

add_cluster(x, c, label = NULL)

get_clusters(x)

remove_cluster(x, i)

## S3 method for class 'DSD_MG'
add_cluster(x, c, label = NULL)

```

### Arguments

dimension	the dimension of the DSD object
...	initial set of <a href="#">MGCs</a>
description	An optional string used by <code>print()</code> to describe the data generator.
x	A DSD_MG object.
c	The cluster that should be added to the DSD_MG object.
label, labels	integer representing the cluster label. NA represents noise. If labels are not specified, then each new cluster gets a new label.
i	The index of the cluster that should be removed from the DSD_MG object.

### Details

This DSD is able to generate complex datasets that are able to evolve over a period of time. Its behavior is determined by as set of [MGCs](#), each representing a moving cluster.

**Author(s)**

Matthew Bolanos

**See Also**[MGC](#) for types of moving clusters.

Other DSD: [DSD\(\)](#), [DSD\\_BarsAndGaussians\(\)](#), [DSD\\_Benchmark\(\)](#), [DSD\\_Cubes\(\)](#), [DSD\\_Gaussians\(\)](#), [DSD\\_Memory\(\)](#), [DSD\\_Mixture\(\)](#), [DSD\\_NULL\(\)](#), [DSD\\_ReadDB\(\)](#), [DSD\\_ReadStream\(\)](#), [DSD\\_Target\(\)](#), [DSD\\_UniformNoise\(\)](#), [DSD\\_mlbenchData\(\)](#), [DSD\\_mlbenchGenerator\(\)](#), [DSF\(\)](#), [animate\\_data\(\)](#), [close\\_stream\(\)](#), [get\\_points\(\)](#), [plot.DSD\(\)](#), [reset\\_stream\(\)](#)

**Examples**

```
### create an empty DSD_MG
stream <- DSD_MG(dim = 2)
stream

### add two clusters
c1 <- MGC_Random(density = 50, center = c(50, 50), parameter = 1)
add_cluster(stream, c1)
stream

c2 <- MGC_Noise(density = 1, range = rbind(c(-20, 120), c(-20, 120)))
add_cluster(stream, c2)
stream

get_clusters(stream)
get_points(stream, n = 5)
plot(stream, xlim = c(-20,120), ylim = c(-20, 120))

if (interactive()) {
  animate_data(stream, n = 5000, xlim = c(-20, 120), ylim = c(-20, 120))
}

### remove cluster 1
remove_cluster(stream, 1)
stream

get_clusters(stream)
plot(stream, xlim = c(-20, 120), ylim = c(-20, 120))

### create a more complicated cluster structure (using 2 clusters with the same
### label to form an L shape)
stream <- DSD_MG(dim = 2,
  MGC_Static(density = 10, center = c(.5, .2), par = c(.4, .2), shape = Shape_Block),
  MGC_Static(density = 10, center = c(.6, .5), par = c(.2, .4), shape = Shape_Block),
  MGC_Static(density = 5, center = c(.39, .53), par = c(.16, .35), shape = Shape_Block),
  MGC_Noise(density = 1, range = rbind(c(0,1), c(0,1))),
  labels = c(1, 1, 2, NA)
)
stream
```

```

plot(stream, xlim = c(0, 1), ylim = c(0, 1))

### simulate the clustering of a splitting cluster
c1 <- MGC_Linear(dim = 2, keyframelist = list(
  keyframe(time = 1, dens = 20, center = c(0,0), param = 10),
  keyframe(time = 50, dens = 10, center = c(50,50), param = 10),
  keyframe(time = 100,dens = 10, center = c(50,100),param = 10)
))

### Note: Second cluster appearch at time=50
c2 <- MGC_Linear(dim = 2, keyframelist = list(
  keyframe(time = 50, dens = 10, center = c(50,50), param = 10),
  keyframe(time = 100,dens = 10, center = c(100,50),param = 10)
))

stream <- DSD_MG(dim = 2, c1, c2)
stream

dbstream <- DSC_DBSTREAM(r = 20, lambda = 0.1)
if (interactive()) {
  purity <- animate_cluster(dbstream, stream, n = 2500, type = "micro",
    xlim = c(-10, 120), ylim = c(-10, 120), measure = "purity", horizon = 100)
}

```

---

DSD\_Mixture

*Mixes Data Points from Several Streams into a Single Stream*


---

## Description

This generator mixes multiple streams given specified probabilities. The streams have to contain the same number of dimensions.

## Usage

```
DSD_Mixture(..., prob = NULL)
```

## Arguments

...	<a href="#">DSD</a> objects.
prob	a numeric vector with the probability for each stream that the next point will be drawn from that stream.

## Value

Returns a `DSD_Mixture` object.(subclass of [DSD\\_R](#), [DSD](#)).

## Author(s)

Michael Hahsler

**See Also**

Other DSD: [DSD\(\)](#), [DSD\\_BarsAndGaussians\(\)](#), [DSD\\_Benchmark\(\)](#), [DSD\\_Cubes\(\)](#), [DSD\\_Gaussians\(\)](#), [DSD\\_MG\(\)](#), [DSD\\_Memory\(\)](#), [DSD\\_NULL\(\)](#), [DSD\\_ReadDB\(\)](#), [DSD\\_ReadStream\(\)](#), [DSD\\_Target\(\)](#), [DSD\\_UniformNoise\(\)](#), [DSD\\_mlbenchData\(\)](#), [DSD\\_mlbenchGenerator\(\)](#), [DSF\(\)](#), [animate\\_data\(\)](#), [close\\_stream\(\)](#), [get\\_points\(\)](#), [plot.DSD\(\)](#), [reset\\_stream\(\)](#)

**Examples**

```
# create data stream with three clusters in 2D
stream1 <- DSD_Gaussians(d = 2, k = 3)
stream2 <- DSD_UniformNoise(d = 2, range = rbind(c(-.5, 1.5), c(-.5, 1.5)))

combinedStream <- DSD_Mixture(stream1, stream2, prob = c(.9, .1))
combinedStream

get_points(combinedStream, n = 20)
plot(combinedStream, n = 200)
```

---

DSD\_mlbenchData

*Stream Interface for Data Sets From mlbench*


---

**Description**

Provides a convenient stream interface for data sets from the mlbench package.

**Usage**

```
DSD_mlbenchData(data = NULL, loop = FALSE, random = FALSE, scale = FALSE)
```

**Arguments**

data	The name of the dataset from mlbench. If missing then a list of all available data sets is shown and returned.
loop	logical; loop or not to loop over the data frame.
random	logical; should the data be used a random order?
scale	logical; apply scaling to the data?

**Details**

The DSD\_mlbenchData class is designed to be a wrapper class for data from the mlbench package. All data is held in memory in either data frame or matrix form. It is served as a stream using the [DSD\\_Memory](#) class. The stream can be reset to position 1 using [reset\\_stream\(\)](#). Call DSD\_mlbenchData with a missing value for data to get a list of all available data sets.

**Value**

Returns a DSD\_mlbenchData object which is also of class [DSD\\_Memory](#).

**Author(s)**

Michael Hahsler and Matthew Bolanos

**See Also**

Other DSD: [DSD\(\)](#), [DSD\\_BarsAndGaussians\(\)](#), [DSD\\_Benchmark\(\)](#), [DSD\\_Cubes\(\)](#), [DSD\\_Gaussians\(\)](#), [DSD\\_MG\(\)](#), [DSD\\_Memory\(\)](#), [DSD\\_Mixture\(\)](#), [DSD\\_NULL\(\)](#), [DSD\\_ReadDB\(\)](#), [DSD\\_ReadStream\(\)](#), [DSD\\_Target\(\)](#), [DSD\\_UniformNoise\(\)](#), [DSD\\_mlbenchGenerator\(\)](#), [DSF\(\)](#), [animate\\_data\(\)](#), [close\\_stream\(\)](#), [get\\_points\(\)](#), [plot.DSD\(\)](#), [reset\\_stream\(\)](#)

**Examples**

```
DSD_mlbenchData()

stream <- DSD_mlbenchData("Shuttle")
stream

get_points(stream, n = 5)

plot(stream, n = 100)
```

---

DSD\_mlbenchGenerator    *mlbench Data Stream Generator*

---

**Description**

A data stream generator class that interfaces data generators found in package mlbench.

**Usage**

```
DSD_mlbenchGenerator(method, ...)
```

**Arguments**

method	The name of the mlbench data generator. If missing then a list of all available generators is shown and returned.
...	Parameters for the mlbench data generator.

**Details**

The `DSD_mlbenchGenerator` class is designed to be a wrapper class for data created by data generators in the mlbench library.

Call `DSD_mlbenchGenerator` with missing `method` to get a list of available methods.

**Value**

Returns a `DSD_mlbenchGenerator` object (subclass of [DSD\\_R](#), [DSD](#))

**Author(s)**

John Forrest

**See Also**

Other DSD: [DSD\(\)](#), [DSD\\_BarsAndGaussians\(\)](#), [DSD\\_Benchmark\(\)](#), [DSD\\_Cubes\(\)](#), [DSD\\_Gaussians\(\)](#), [DSD\\_MG\(\)](#), [DSD\\_Memory\(\)](#), [DSD\\_Mixture\(\)](#), [DSD\\_NULL\(\)](#), [DSD\\_ReadDB\(\)](#), [DSD\\_ReadStream\(\)](#), [DSD\\_Target\(\)](#), [DSD\\_UniformNoise\(\)](#), [DSD\\_mlbenchData\(\)](#), [DSF\(\)](#), [animate\\_data\(\)](#), [close\\_stream\(\)](#), [get\\_points\(\)](#), [plot.DSD\(\)](#), [reset\\_stream\(\)](#)

**Examples**

```
DSD_mlbenchGenerator()

stream <- DSD_mlbenchGenerator(method = "cassini")
stream

get_points(stream, n = 5)

plot(stream, n = 500)
```

---

DSD\_NULL

*Placeholder for a DSD Stream*

---

**Description**

Placeholder for a [DSD](#). DSD\_NULL does not produce points and creates an error for [get\\_points\(\)](#).

**Usage**

```
DSD_NULL()
```

**Value**

Returns a DSD\_NULL object (subclass of [DSD](#)).

**Author(s)**

Michael Hahsler

**See Also**

Other DSD: [DSD\(\)](#), [DSD\\_BarsAndGaussians\(\)](#), [DSD\\_Benchmark\(\)](#), [DSD\\_Cubes\(\)](#), [DSD\\_Gaussians\(\)](#), [DSD\\_MG\(\)](#), [DSD\\_Memory\(\)](#), [DSD\\_Mixture\(\)](#), [DSD\\_ReadDB\(\)](#), [DSD\\_ReadStream\(\)](#), [DSD\\_Target\(\)](#), [DSD\\_UniformNoise\(\)](#), [DSD\\_mlbenchData\(\)](#), [DSD\\_mlbenchGenerator\(\)](#), [DSF\(\)](#), [animate\\_data\(\)](#), [close\\_stream\(\)](#), [get\\_points\(\)](#), [plot.DSD\(\)](#), [reset\\_stream\(\)](#)

**Examples**

```

nullstream <- DSD_NULL()
nullstream

## This will produce an error
## Not run:
get_points(nullstream)
## End(Not run)

```

---

DSD\_ReadDB

*Read a Data Stream from an open DB Query*


---

**Description**

A DSD class that reads a data stream from an open DB result set from a relational database with using R's data base interface (DBI).

**Usage**

```

DSD_ReadDB(
  result,
  k = NA,
  outofpoints = c("warn", "ignore", "stop"),
  description = NULL
)

## S3 method for class 'DSD_ReadDB'
close_stream(dsd, disconnect = TRUE, ...)

```

**Arguments**

result	An open DBI result set.
k	Number of true clusters, if known.
outofpoints	Action taken if less than n data points are available. The default is to return the available data points with a warning. Other supported actions are: <ul style="list-style-type: none"> <li>• warn: return the available points (maybe an empty data.frame) with a warning.</li> <li>• ignore: silently return the available points.</li> <li>• stop: stop with an error.</li> </ul>
description	a character string describing the data.
dsd	a stream.
disconnect	logical; disconnect from the database?
...	further arguments.

## Details

This class provides a streaming interface for result sets from a data base with via [DBI::DBI](#). You need to connect to the data base and submit a SQL query using [DBI::dbGetQuery\(\)](#) to obtain a result set. Make sure that your query only includes the columns that should be included in the stream (including class and outlier marking columns).

### Closing and resetting the stream

Do not forget to clear the result set and disconnect from the data base connection. [close\\_stream\(\)](#) clears the query result with [DBI::dbClearResult\(\)](#) and the disconnects from the database with [DBI::dbDisconnect\(\)](#). Disconnecting can be prevented by calling [close\\_stream\(\)](#) with `disconnect = FALSE`.

[reset\\_stream\(\)](#) is not available for this type of stream.

### Additional information

If additional information is available (e.g., class information), then the SQL statement needs to make sure that the columns have the appropriate name starting with `.`. See Examples section below.

## Value

An object of class `DSD_ReadDB` (subclass of [DSD\\_R](#), [DSD](#)).

## Author(s)

Michael Hahsler

## See Also

[DBI::dbGetQuery\(\)](#)

Other DSD: [DSD\(\)](#), [DSD\\_BarsAndGaussians\(\)](#), [DSD\\_Benchmark\(\)](#), [DSD\\_Cubes\(\)](#), [DSD\\_Gaussians\(\)](#), [DSD\\_MG\(\)](#), [DSD\\_Memory\(\)](#), [DSD\\_Mixture\(\)](#), [DSD\\_NULL\(\)](#), [DSD\\_ReadStream\(\)](#), [DSD\\_Target\(\)](#), [DSD\\_UniformNoise\(\)](#), [DSD\\_mlbenchData\(\)](#), [DSD\\_mlbenchGenerator\(\)](#), [DSF\(\)](#), [animate\\_data\(\)](#), [close\\_stream\(\)](#), [get\\_points\(\)](#), [plot.DSD\(\)](#), [reset\\_stream\(\)](#)

## Examples

```
### create a data base with a table with 3 Gaussians

library("DBI")
con <- dbConnect(RSQLite::SQLite(), ":memory:")

points <- get_points(DSD_Gaussians(k = 3, d = 2), n = 110)
head(points)

dbWriteTable(con, "Gaussians", points)

### prepare a query result set. Make sure that the additional information
### column starts with .
res <- dbSendQuery(con, "SELECT X1, X2, `.class` AS '.class' FROM Gaussians")
res
```



```
### create a stream interface to the result set
stream <- DSD_ReadDB(res, k = 3)
stream

### get points
get_points(stream, n = 5)

plot(stream, n = 100)

### close stream clears the query and disconnects the database
close_stream(stream)
```

---

DSD_ReadStream	<i>Read a Data Stream from a File or a Connection</i>
----------------	---

---

### Description

A DSD class that reads a data stream (text format) from a file or any R connection.

### Usage

```
DSD_ReadStream(  
  file,  
  k = NA,  
  take = NULL,  
  sep = ",",  
  header = FALSE,  
  skip = 0,  
  col.names = NULL,  
  colClasses = NA,  
  outofpoints = c("warn", "ignore", "stop"),  
  ...  
)
```

```
DSD_ReadCSV(  
  file,  
  k = NA,  
  take = NULL,  
  sep = ",",  
  header = FALSE,  
  skip = 0,  
  col.names = NULL,  
  colClasses = NA,  
  outofpoints = c("warn", "ignore", "stop"),  
  ...  
)
```

```
## S3 method for class 'DSD_ReadStream'
close_stream(dsd, ...)
```

```
## S3 method for class 'DSD_ReadCSV'
close_stream(dsd, ...)
```

## Arguments

<code>file</code>	A file/URL or an open connection.
<code>k</code>	Number of true clusters, if known.
<code>take</code>	indices of columns to extract from the file.
<code>sep</code>	The character string that separates dimensions in data points in the stream.
<code>header</code>	Does the first line contain variable names?
<code>skip</code>	the number of lines of the data file to skip before beginning to read data.
<code>col.names</code>	A vector of optional names for the variables. The default is to use "V" followed by the column number. Additional information (e.g., class labels) need to have names starting with ..
<code>colClasses</code>	A vector of classes to be assumed for the columns passed on to <code>read.table()</code> .
<code>outofpoints</code>	Action taken if less than n data points are available. The default is to return the available data points with a warning. Other supported actions are: <ul style="list-style-type: none"> <li>• <code>warn</code>: return the available points (maybe an empty data.frame) with a warning.</li> <li>• <code>ignore</code>: silently return the available points.</li> <li>• <code>stop</code>: stop with an error.</li> </ul>
<code>...</code>	Further arguments are passed on to <code>read.table()</code> . This can for example be used for encoding, quotes, etc.
<code>dsd</code>	A object of class <code>DSD_ReadCSV</code> .

## Details

`DSD_ReadStream` uses `readLines()` and `read.table()` to read data from an R connection line-by-line and convert it into a data.frame. The connection is responsible for maintaining where the stream is currently being read from. In general, the connections will consist of files stored on disk but have many other possibilities (see [connection](#)).

The implementation tries to gracefully deal with slightly corrupted data by dropping points with inconsistent reading and producing a warning. However, this might not always be possible resulting in an error instead.

### Column names

If the file has column headers in the first line, then they can be used by setting `header = TRUE`. Alternatively, column names can be set using `col.names` or a named vector for `take`. If no column names are specified then default names will be created.

Columns with names that start with `.` are considered information columns and are ignored by DSTs. See `get_points()` for details.

Other information columns are used by various functions.

**Reading the whole stream** By using `n = -1` in `get_points()`, the whole stream is returned.

### Resetting and closing a stream

The position in the file can be reset to the beginning or another position using `reset_stream()`. This fails if the underlying connection is not seekable (see [connection](#)).

DSD\_ReadStream maintains an open connection to the stream and needs to be closed using `close_stream()`.

DSD\_ReadCSV reads a stream from a comma-separated values file.

### Value

An object of class DSD\_ReadCSV (subclass of [DSD\\_R](#), [DSD](#)).

### Author(s)

Michael Hahsler

### See Also

[readLines\(\)](#), [read.table\(\)](#).

Other DSD: [DSD\(\)](#), [DSD\\_BarsAndGaussians\(\)](#), [DSD\\_Benchmark\(\)](#), [DSD\\_Cubes\(\)](#), [DSD\\_Gaussians\(\)](#), [DSD\\_MG\(\)](#), [DSD\\_Memory\(\)](#), [DSD\\_Mixture\(\)](#), [DSD\\_NULL\(\)](#), [DSD\\_ReadDB\(\)](#), [DSD\\_Target\(\)](#), [DSD\\_UniformNoise\(\)](#), [DSD\\_mlbenchData\(\)](#), [DSD\\_mlbenchGenerator\(\)](#), [DSF\(\)](#), [animate\\_data\(\)](#), [close\\_stream\(\)](#), [get\\_points\(\)](#), [plot.DSD\(\)](#), [reset\\_stream\(\)](#)

### Examples

```
# Example 1: creating data and writing it to disk
stream <- DSD_Gaussians(k = 3, d = 2)
write_stream(stream, "data.txt", n = 100, info = TRUE, header = TRUE)
readLines("data.txt", n = 5)

# reading the same data back
stream2 <- DSD_ReadStream("data.txt", header = TRUE)
stream2

# get points
get_points(stream2, n = 5)
plot(stream2, n = 20)

# clean up
close_stream(stream2)
file.remove("data.txt")

# Example 2: Read part of the kddcup1999 data (take only cont. variables)
# col 42 is the class variable
file <- system.file("examples", "kddcup10000.data.gz", package = "stream")
stream <- DSD_ReadCSV(gzfile(file),
  take = c(1, 5, 6, 8:11, 13:20, 23:41, .class = 42), k = 7)
stream

get_points(stream, 5)
```

```
# plot 100 points (projected on the first two principal components)
plot(stream, n = 100, method = "pca")

close_stream(stream)
```

---

DSD\_ScaleStream      *Deprecated DSD\_ScaleStream*

---

### Description

Deprecated! Use [DSF\\_Scale](#) instead.

### Usage

```
DSD_ScaleStream(...)
```

### Arguments

...                    arguments are passed on to [DSF\\_Scale\(\)](#).

### Value

Produces a warning that DSD\_ScaleStream is deprecated and returns an object of class [DSF\\_Scale](#) (subclass of [DSF](#) and [DSD](#)).

---

DSD\_Target              *Target Data Stream Generator*

---

### Description

A data stream generator that generates a data stream in the shape of a target. It has a single Gaussian cluster in the center and a ring that surrounds it.

### Usage

```
DSD_Target(
  center_sd = 0.05,
  center_weight = 0.5,
  ring_r = 0.2,
  ring_sd = 0.02,
  noise = 0
)
```

**Arguments**

center_sd	standard deviation of center
center_weight	proportion of points in center
ring_r	average ring radius
ring_sd	standard deviation of ring radius
noise	proportion of noise

**Details**

This DSD will produce a singular Gaussian cluster in the center with a ring around it.

**Value**

Returns a DSD\_Target object.

**Author(s)**

Michael Hahsler

**See Also**

Other DSD: [DSD\(\)](#), [DSD\\_BarsAndGaussians\(\)](#), [DSD\\_Benchmark\(\)](#), [DSD\\_Cubes\(\)](#), [DSD\\_Gaussians\(\)](#), [DSD\\_MG\(\)](#), [DSD\\_Memory\(\)](#), [DSD\\_Mixture\(\)](#), [DSD\\_NULL\(\)](#), [DSD\\_ReadDB\(\)](#), [DSD\\_ReadStream\(\)](#), [DSD\\_UniformNoise\(\)](#), [DSD\\_mlbenchData\(\)](#), [DSD\\_mlbenchGenerator\(\)](#), [DSF\(\)](#), [animate\\_data\(\)](#), [close\\_stream\(\)](#), [get\\_points\(\)](#), [plot.DSD\(\)](#), [reset\\_stream\(\)](#)

**Examples**

```
# create data stream with three clusters in 2D
stream <- DSD_Target()

get_points(stream, n = 5)

plot(stream)
```

---

DSD\_UniformNoise

*Uniform Noise Data Stream Generator*

---

**Description**

This generator produces uniform noise in a d-dimensional unit (hyper) cube.

**Usage**

```
DSD_UniformNoise(d = 2, range = NULL)
```

**Arguments**

`d` Determines the number of dimensions.

`range` A matrix with two columns and `d` rows giving the minimum and maximum for each dimension. Defaults to the range of  $[0, 1]$ .

**Value**

Returns a `DSD_UniformNoise` object.(subclass of `DSD_R`, `DSD`).

**Author(s)**

Michael Hahsler

**See Also**

Other DSD: [DSD\(\)](#), [DSD\\_BarsAndGaussians\(\)](#), [DSD\\_Benchmark\(\)](#), [DSD\\_Cubes\(\)](#), [DSD\\_Gaussians\(\)](#), [DSD\\_MG\(\)](#), [DSD\\_Memory\(\)](#), [DSD\\_Mixture\(\)](#), [DSD\\_NULL\(\)](#), [DSD\\_ReadDB\(\)](#), [DSD\\_ReadStream\(\)](#), [DSD\\_Target\(\)](#), [DSD\\_mlbenchData\(\)](#), [DSD\\_mlbenchGenerator\(\)](#), [DSF\(\)](#), [animate\\_data\(\)](#), [close\\_stream\(\)](#), [get\\_points\(\)](#), [plot.DSD\(\)](#), [reset\\_stream\(\)](#)

**Examples**

```
# create data stream with three clusters in 2D
stream <- DSD_UniformNoise(d = 2)
get_points(stream, n = 5)
plot(stream, n = 100)

# specify a different range for each dimension
stream <- DSD_UniformNoise(d = 3,
  range = rbind(c(0, 1), c(0, 10), c(0, 5)))
plot(stream, n = 100)
```

---

DSF

*Data Stream Filter Base Classes*


---

**Description**

Abstract base classes for all data stream filter (DSF) classes. Data stream filters transform a data stream ([DSD](#)).

**Usage**

```
DSF(...)
```

```
## S3 method for class 'DSF'
reset_stream(dsd, pos = 1)
```

```
## S3 method for class 'DSF'
```

```

get_points(x, n = 1L, info = TRUE, ...)

## S3 method for class 'DSF'
update(object, dsd = NULL, n = 1L, return = "data", ...)

## S3 method for class 'DSF'
close_stream(dsd, ...)

```

### Arguments

...	Further arguments passed on.
dsd	a stream object of class <a href="#">DSD</a> .
pos	position in the stream.
x, object	a DSF object.
n	number of points to get/use for the update.
info	return additional columns with information about the data point (e.g., a known cluster assignment).
return	a character string indicating what update returns. The only value is currently "data" to return the transformed data. possible values depend on the DST.

### Details

The DSF class cannot be instantiated, but it serve as a base class from which other DSF classes inherit.

Data stream filters transform a [DSD](#) data stream. DSF can be used in two ways.

1. **DSD Adapter:** When a data stream (dsd) is specified in the constructor, then the DSF acts as an adapter for a connected data stream. The DSF implementations inherit the interface from [DSD](#) and provide:
  - [get\\_points\(\)](#) get the transformed points.
  - [reset\\_stream\(\)](#) reset the underlying stream
  - [close\\_stream\(\)](#) close the underlying stream
2. **Stream Transformer:** When no data stream (dsd) is specified in the constructor, then the DSF acts like a [DST](#) data stream task and provides:
  - [update\(\)](#) to transform the points from a specified DSD. It is convenient to use the pipe ([magrittr::%>%](#)) to apply one or more filters to data streams (see Examples section).

### Methods (by generic)

- [reset\\_stream\(DSF\)](#): reset the attached stream if reset is supported.
- [get\\_points\(DSF\)](#): DSD-like interface to get points if the DSF was created with an attached stream.
- [update\(DSF\)](#): updates with data and returns the filtered data.
- [close\\_stream\(DSF\)](#): close the attached stream if close is supported.

**Author(s)**

Michael Hahsler

**See Also**

Other DSF: [DSF\\_Convolve\(\)](#), [DSF\\_Downsampling\(\)](#), [DSF\\_ExponentialMA\(\)](#), [DSF\\_FeatureSelection\(\)](#), [DSF\\_Func\(\)](#), [DSF\\_Scale\(\)](#), [DSF\\_dplyr\(\)](#)

Other DSD: [DSD\(\)](#), [DSD\\_BarsAndGaussians\(\)](#), [DSD\\_Benchmark\(\)](#), [DSD\\_Cubes\(\)](#), [DSD\\_Gaussians\(\)](#), [DSD\\_MG\(\)](#), [DSD\\_Memory\(\)](#), [DSD\\_Mixture\(\)](#), [DSD\\_NULL\(\)](#), [DSD\\_ReadDB\(\)](#), [DSD\\_ReadStream\(\)](#), [DSD\\_Target\(\)](#), [DSD\\_UniformNoise\(\)](#), [DSD\\_mlbenchData\(\)](#), [DSD\\_mlbenchGenerator\(\)](#), [animate\\_data\(\)](#), [close\\_stream\(\)](#), [get\\_points\(\)](#), [plot.DSD\(\)](#), [reset\\_stream\(\)](#)

**Examples**

```
DSF()

# Example 1: Use as a DSD adapter
stream <- DSD_Gaussians(k = 3, d = 2) %>%
  DSF_Func(func = function(x) cbind(x, Xsum = x$X1 + x$X2))
stream

get_points(stream, n = 5)

# Example 2: Use as a stream transformer
trans <- DSF_Func(func = function(x) cbind(x, Xsum = x$X1 + x$X2))
trans

update(trans, stream, n = 5)

# Example 3: Use as a DST preprocessor
clusterer <- DSF_Func(func = function(x) cbind(x, X1_squared = x$X1^2)) %>%
  DST_Runner(DSC_Kmeans(k = 3))
clusterer

update(clusterer, stream, n = 100)

# Example 5: Specify a complete pipeline DSD -> DSF -> DST
pipeline <- DSD_Gaussians(k = 3, d = 2) %>%
  DSF_Func(func = function(x) cbind(x, X1_squared = x$X1^2)) %>%
  DST_Runner(DSC_Kmeans(k = 3))
pipeline

update(pipeline, n = 100)
plot(pipeline$dst)
```



**Description**

Abstract class for frequent pattern mining algorithms for data streams. Currently, **stream** does not implement frequent pattern mining algorithms.

**Usage**

```
DSFP(...)
```

**Arguments**

```
...          Further arguments.
```

**Author(s)**

Michael Hahsler

**See Also**

[DST](#)

**Examples**

```
DSFP()
```

---

DSF\_Convolve

*Apply a Filter to a Data Stream*

---

**Description**

Applies a filter (i.e., a convolution with a filter kernel) to a data stream.

**Usage**

```
DSF_Convolve(  
  dsd = NULL,  
  dim = NULL,  
  kernel = NULL,  
  pre = NULL,  
  post = NULL,  
  na.rm = FALSE,  
  replace = TRUE,  
  name = NULL  
)  
  
filter_MA(width)  
  
filter_Hamming(width)
```

```
filter_diff(lag)
```

```
filter_Sinc(fc, fs, width = NULL, bw = NULL)
```

```
pow2(x)
```

### Arguments

<code>dsd</code>	A object of class <a href="#">DSD</a> .
<code>dim</code>	columns to which the filter should be applied. Default is all columns.
<code>kernel</code>	filter kernel as a numeric vector of weights.
<code>pre, post</code>	functions to be applied before and after the convolution.
<code>na.rm</code>	logical; should NAs be ignored?
<code>replace</code>	logical; should the column be replaced or a column with the convolved column added?
<code>name</code>	character; the new column will be name with the old column name + <code>_</code> + name.
<code>width</code>	filter width.
<code>lag</code>	an integer indicating which time lag to use.
<code>fc</code>	cutoff frequency.
<code>fs</code>	sampling frequency.
<code>bw</code>	transition bandwidth.
<code>x</code>	values to be squared.

### Details

A filter kernel is a vector with kernel weights. A few filter are provided.

- `filter_MA(width)` creates a moving average.
- `filter_diff(lag)` calculates lagged differences. Note that `na.rm = TRUE` will lead to artifacts and should not be used.
- `filter_Hamming(width)` creates a Hamming window.
- `filter_Sinc(fc, fs, width, bw)` creates a windowed-sinc filter. One of `width` (filter length) or `bw` (transition bandwidth) can be used to control the filter roll-off. The relationship is  $width = 4/bw$ . See Chapter 16 in Smith (1997).

`pre` and `post` are functions that are called before and after the convolution. For example, to calculate RMS, you can use `pre = pow2` and `post = sqrt`. `pow2()` is a convenience function.

### Value

An object of class `DSF_Convolve` (subclass of [DSF](#) and [DSD](#)).

### Author(s)

Michael Hahsler

## References

Steven W. Smith, The Scientist and Engineer's Guide to Digital Signal Processing, California Technical Pub; 1st edition (January 1, 1997). ISBN 0966017633, URL: <https://www.dspguide.com/>

## See Also

`stats::filter` provides non-streaming convolution.

Other DSF: `DSF()`, `DSF_Downsamp()`, `DSF_ExponentialMA()`, `DSF_FeatureSelection()`, `DSF_Func()`, `DSF_Scale()`, `DSF_dplyr()`

## Examples

```
data(presidents)

## Example 1: Create a data stream with three copies of president approval ratings.
## We will use several convolutions.
stream <- data.frame(
  approval_orig = presidents,
  approval_MA = presidents,
  approval_diff1 = presidents,
  .time = time(presidents)) %>%
  DSD_Memory()

plot(stream, dim = 1, n = 120, method = "ts")

## apply a moving average filter to dimension 1 (using the column name) and diff to dimension 3
filteredStream <- stream %>%
  DSF_Convolve(kernel = filter_MA(5), dim = "approval_orig", na.rm = TRUE) %>%
  DSF_Convolve(kernel = filter_diff(1), dim = 3)
filteredStream

## resetting the filtered stream also resets the original stream
reset_stream(filteredStream)
ps <- get_points(filteredStream, n = 120)
head(ps)

year <- ps[[".time"]]
approval <- remove_info(ps)
matplot(year, approval, type = "l", ylim = c(-20, 100))
legend("topright", colnames(approval), col = 1:3, lty = 1:3, bty = "n")

## Example 2: Create a stream with a constant sine wave and apply
## a moving average, an RMS envelope and a differences
stream <- DSD_Memory(data.frame(y = sin(seq(0, 2 * pi - (2 * pi / 100) ,
  length.out = 100))), loop = TRUE)
plot(stream, n = 200, method = "ts")

filteredStream <- stream %>%
  DSF_Convolve(kernel = filter_MA(100), dim = 1,
    replace = FALSE, name = "MA") %>%
  DSF_Convolve(kernel = filter_MA(100), pre = pow2, post = sqrt, dim = 1,
```

```

        replace = FALSE, name = "RMS") %>%
    DSF_Convolve(kernel = filter_diff(1), dim = 1,
        replace = FALSE, name = "diff1")
    filteredStream

ps <- get_points(filteredStream, n = 500)
head(ps)

matplot(ps, type = "l")
legend("topright", colnames(ps), col = 1:4, lty = 1:4)

## Note that MA and RMS use a window of length 200 and are missing at the
## beginning of the stream the window is full.

## Filters: look at different filters
filter_MA(5)
filter_diff(1)
plot(filter_Hamming(20), type = "h")
plot(filter_Sinc(10, 100, width = 20), type = "h")

```

---

DSF\_Downsampler

*Downsample a Data Stream*


---

## Description

Creates a new stream that reduces the frequency of a given stream by a given factor.

## Usage

```
DSF_Downsampler(dsd = NULL, factor = 1L)
```

## Arguments

<code>dsd</code>	The input stream as an <a href="#">DSD</a> object.
<code>factor</code>	the downsampling factor.

## Value

An object of class `DSF_Downsampler` (subclass of [DSF](#) and [DSD](#)).

## Author(s)

Michael Hahsler

## See Also

Other DSF: [DSF\(\)](#), [DSF\\_Convolve\(\)](#), [DSF\\_ExponentialMA\(\)](#), [DSF\\_FeatureSelection\(\)](#), [DSF\\_Func\(\)](#), [DSF\\_Scale\(\)](#), [DSF\\_dplyr\(\)](#)

**Examples**

```

# Simple downsampling example
stream <- DSD_Memory(data.frame(rownum = seq(100))) %>% DSD_Downsampling(factor = 10)
stream

get_points(stream, n = 2)
get_points(stream, n = 1)
get_points(stream, n = 5)

# DSD_Memory supports getting the remaining points using n = -1
get_points(stream, n = -1)

# Downsample a time series
data(presidents)

stream <- data.frame(
  presidents,
  .time = time(presidents)) %>%
  DSD_Memory()

plot(stream, dim = 1, n = 120, method = "ts")

# downsample by taking only every 3rd data point (quarters)
downsampledStream <- stream %>% DSD_Downsampling(factor = 3)

reset_stream(downsampledStream)
plot(downsampledStream, dim = 1, n = 40, method = "ts")

```

DSF\_dplyr

*Apply a dplyr Transformation to a Data Stream***Description**

Applies dplyr transformations to a data stream.

**Usage**

```
DSF_dplyr(dsd = NULL, func = NULL, info = FALSE)
```

**Arguments**

dsd	A object of class <a href="#">DSD</a> .
func	a dplyr expression.
info	logical; does the function also receive and modify the info columns?

## Details

**dplyr** needs to be installed and loaded with `library(dplyr)` before `DSF_dplyr` can be used.

Since streams are processed one point or block at a time, only `dplyr::dplyr` operations that work on individual rows are allowed on streams. Examples are:

- `dplyr::select()`
- `dplyr::mutate()`
- `dplyr::rename()`
- `dplyr::transmute()`
- `dplyr::filter()`

Summary functions can be used, but will only be applied to the requested part of the stream of length `n`.

`DSF_dplyr()` calls the function using points `%>% <func>` and multiple `dplyr` functions can be applied by using `%>%` between them.

## Value

An object of class `DSF_dplyr` (subclass of `DSF` and `DSD`).

## Author(s)

Michael Hahsler

## See Also

Other DSF: `DSF()`, `DSF_Convolve()`, `DSF_Downsampling()`, `DSF_ExponentialMA()`, `DSF_FeatureSelection()`, `DSF_Func()`, `DSF_Scale()`

## Examples

```
library(dplyr)

stream <- DSD_Gaussians(k = 3, d = 3)
plot(stream, xlim = c(0, 1), ylim = c(0, 1))

# 1. Select only columns X1 and X2
# 2. filter points by X1 > .5 (Note that the info columns also need to be filtered!)
# 3. Add a sum columns

stream2 <- stream %>%
  DSF_dplyr(select(X1, X2)) %>%
  DSF_dplyr(filter(X1 > .5), info = TRUE) %>%
  DSF_dplyr(mutate(Xsum = X1 + X2))
stream2

# Note: you get fewer points because of the filter operation.
get_points(stream2, n = 10)
```

```
plot(stream2, xlim = c(0, 1), ylim = c(0, 1))
```

---

DSF\_ExponentialMA      *Exponential Moving Average over a Data Stream*

---

### Description

Applies an exponential moving average to components of a data stream.

### Usage

```
DSF_ExponentialMA(dsd = NULL, dim = NULL, alpha = 0.5)
```

### Arguments

dsd	The input stream as an <a href="#">DSD</a> object.
dim	columns to which the filter should be applied. Default is all columns.
alpha	smoothing coefficient in $[0, 1]$ . Larger means discounting older observations faster.

### Details

The exponential moving average is calculated by:

$$S_t = \alpha Y_t + (1 - \alpha) S_{t-1}$$

with  $S_0 = Y_0$ .

### Value

An object of class `DSF_ExponentialMA` (subclass of [DSF](#) and [DSD](#)).

### Author(s)

Michael Hahsler

### See Also

Other DSF: [DSF\(\)](#), [DSF\\_Convolve\(\)](#), [DSF\\_Downsampling\(\)](#), [DSF\\_FeatureSelection\(\)](#), [DSF\\_Func\(\)](#), [DSF\\_Scale\(\)](#), [DSF\\_dplyr\(\)](#)

## Examples

```
# Smooth a time series
data(presidents)

stream <- data.frame(
  presidents,
  .time = time(presidents)) %>%
  DSD_Memory()

plot(stream, dim = 1, n = 120, method = "ts", main = "Original")

smoothStream <- stream %>% DSF_ExponentialMA(alpha = .7)
smoothStream

reset_stream(smoothStream)
plot(smoothStream, dim = 1, n = 120, method = "ts", main = "With ExponentialMA(.7)")
```

---

DSF\_FeatureSelection *Select Features for a Data Stream*

---

## Description

Select features from a data stream given a list of features.

## Usage

```
DSF_FeatureSelection(dsd = NULL, features)
```

## Arguments

<code>dsd</code>	A object of class <a href="#">DSD</a> that will be scaled.
<code>features</code>	a character vector with feature (column) names or the numeric index of the selected features. All other features will be removed. Note special info columns starting with <code>.</code> are not features.

## Value

An object of class `DSF_FeatureSelection` (subclass of [DSF](#) and [DSD](#)).

## Author(s)

Michael Hahsler

## See Also

Other DSF: [DSF\(\)](#), [DSF\\_Convolve\(\)](#), [DSF\\_Downsampling\(\)](#), [DSF\\_ExponentialMA\(\)](#), [DSF\\_Func\(\)](#), [DSF\\_Scale\(\)](#), [DSF\\_dplyr\(\)](#)



**Examples**

```
stream <- DSD_Gaussians(k = 3, d = 3)
get_points(stream, 3)

stream_2features <- DSF_FeatureSelection(stream, features = c("X1", "X3"))
stream_2features

get_points(stream_2features, n = 3)
```

---

**DSF\_Func***Apply a Function to Transformation to a Data Stream*

---

**Description**

Applies an R function to transform to a data stream.

**Usage**

```
DSF_Func(dsd = NULL, func, ..., info = FALSE)
```

**Arguments**

<code>dsd</code>	A object of class <a href="#">DSD</a> .
<code>func</code>	a function that takes a data.frame as the first argument and returns the transformed data.frame.
<code>...</code>	further arguments are passed on to the function specified in <code>func</code> .
<code>info</code>	logical; does the function also receive and modify the info columns?

**Details**

The function's first argument needs to be a data.frame representing points of the data stream. The function will be called as `ps %>% your_function()`, where `ps` is the data.frame with some points obtained using [get\\_points\(\)](#) on the data stream source.

**Value**

An object of class `DSF_Func` (subclass of [DSF](#) and [DSD](#)).

**Author(s)**

Michael Hahsler

**See Also**

Other DSF: [DSF\(\)](#), [DSF\\_Convolve\(\)](#), [DSF\\_Downsampling\(\)](#), [DSF\\_ExponentialMA\(\)](#), [DSF\\_FeatureSelection\(\)](#), [DSF\\_Scale\(\)](#), [DSF\\_dplyr\(\)](#)

**Examples**

```

stream <- DSD_Gaussians(k = 3, d = 3)
get_points(stream, n = 5)

## Example 1: rename the columns
rename <- function(x, names) {
  colnames(x) <- names
  x
}

# By default, the info columns starting with . are not affected.
stream2 <- stream %>% DSF_Func(rename, names = c("A", "B", "C"))
stream2
get_points(stream2, n = 5)

## Example 2: add a sum columns
stream3 <- stream2 %>% DSF_Func(function(x) {
  x$sum = rowSums(x)
  x
})
stream3
get_points(stream3, n = 5)

## Example 3: Project the stream on its first 2 PCs (using a sample)
pr <- princomp(get_points(stream, n = 100, info = FALSE))
pca_trans <- function(x) predict(pr, x[, c("X1", "X2", "X3")])[, 1:2, drop = FALSE]
pca_trans(get_points(stream, n = 3, info = FALSE))

stream4 <- stream %>% DSF_Func(pca_trans)
stream4

get_points(stream4, n = 3)
plot(stream4)

## Example 4: Change a class labels using info = TRUE. We redefine class 3 as noise (NA)
stream5 <- stream %>% DSF_Func(
  function(x) { x[['.class']][x[['.class']] == 3] <- NA; x },
  info = TRUE)
stream5

get_points(stream5, n = 5)
plot(stream5)

```

**Description**

Make an unscaled data stream into a scaled data stream.

**Usage**

```
DSF_Scale(dsd = NULL, dim = NULL, center = TRUE, scale = TRUE, n = 100L)
```

**Arguments**

<code>dsd</code>	A object of class <a href="#">DSD</a> that will be scaled.
<code>dim</code>	integer vector or names of dimensions that should be scaled? Default is all.
<code>center, scale</code>	logical or a numeric vector of length equal to the number of columns (selected with <code>dim</code> ) used for centering/scaling (see function <a href="#">scale</a> ).
<code>n</code>	The number of points used by <code>scale_stream()</code> to creating the centering/scaling

**Details**

If `center` and `scale` are not vectors with scaling factors, then `scale_stream()` estimates the values for centering and scaling (see [scale](#) in **base**) using `n` points from the stream and the stream is reset if `reset = TRUE` and the [DSD](#) object supports resetting.

**Value**

An object of class `DSF_Scale` (subclass of [DSF](#) and [DSD](#)).

**Deprecated**

`DSD_ScaleStream` is deprecated. Use `DSF_Scale` instead.

**Author(s)**

Michael Hahsler

**See Also**

[scale](#) in **base**

Other DSF: [DSF\(\)](#), [DSF\\_Convolve\(\)](#), [DSF\\_Downsampling\(\)](#), [DSF\\_ExponentialMA\(\)](#), [DSF\\_FeatureSelection\(\)](#), [DSF\\_Func\(\)](#), [DSF\\_dplyr\(\)](#)

**Examples**

```
stream <- DSD_Gaussians(k = 3, d = 2)
get_points(stream, 3)

# scale with manually calculated scaling factors
points <- get_points(stream, n = 100, info = FALSE)
center <- colMeans(points)
scale <- apply(points, MARGIN = 2, sd)

scaledStream <- stream %>% DSF_Scale(dim = c(1L, 2L), center = center, scale = scale)
colMeans(get_points(scaledStream, n = 100, info = FALSE))
apply(get_points(scaledStream, n = 100, info = FALSE), MARGIN = 2, sd)
```

```
# let DSF_Scale calculate the scaling factors from the first n points of the stream
scaledStream <- stream %>% DSF_Scale(dim = c(1L, 2L), n = 100)
colMeans(get_points(scaledStream, n = 100, info = FALSE))
apply(get_points(scaledStream, n = 100, info = FALSE), MARGIN = 2, sd)

## scale only X2
scaledStream <- stream %>% DSF_Scale(dim = "X2", n = 100)
colMeans(get_points(scaledStream, n = 100, info = FALSE))
apply(get_points(scaledStream, n = 100, info = FALSE), MARGIN = 2, sd)
```

---

DSOutlier

*Abstract Class for Data Stream Outlier Detectors*


---

## Description

The abstract class for all data stream outlier detectors. Cannot be instantiated. Some [DSC](#) implementations also implement outlier/noise detection.

## Usage

```
DSOutlier(...)
```

## Arguments

... further arguments.

## Details

`plot()` has an extra logical argument to specify if outliers should be plotted as red crosses.

## Author(s)

Michael Hahsler

## See Also

Other DST: [DSAggregate\(\)](#), [DSC\(\)](#), [DSCClassifier\(\)](#), [DSRegressor\(\)](#), [DST\(\)](#), [DST\\_SlidingWindow\(\)](#), [DST\\_WriteStream\(\)](#), [evaluate](#), [predict\(\)](#), [stream\\_pipeline](#), [update\(\)](#)

Other DSOutlier: [DSC\\_DBSTREAM\(\)](#), [DSC\\_DStream\(\)](#)

## Examples

```
DSOutlier()

#' @examples
set.seed(1000)
stream <- DSD_Gaussians(k = 3, d = 2, noise = 0.1, noise_separation = 5)

outlier_detector <- DSOutlier_DBSTREAM(r = .05, outlier_multiplier = 2)
```

```

update(outlier_detector, stream, 500)
outlier_detector

points <- get_points(stream, 20)
points

# Outliers are predicted as class NA
predict(outlier_detector, points)

# Plot new points from the stream. Predicted outliers are marked with a red x.
plot(outlier_detector, stream)

evaluate_static(outlier_detector, stream, measure =
  c("noiseActual", "noisePredicted", "noisePrecision", "outlierJaccard"))

# use a different detector
outlier_detector2 <- DSOutlier_DStream(gridsize = .05, Cl = 0.5, outlier_multiplier = 2)
update(outlier_detector2, stream, 500)
plot(outlier_detector2, stream)

evaluate_static(outlier_detector2, stream, measure =
  c("noiseActual", "noisePredicted", "noisePrecision", "outlierJaccard"))

```

---

DSRegressor

*Abstract Class for Data Stream Regressors*


---

## Description

Abstract class for data stream regressors. More implementations can be found in package **stream-MOA**.

## Usage

```
DSRegressor(...)
```

## Arguments

... Further arguments.

## Author(s)

Michael Hahsler

## See Also

Other DST: [DSAggregate\(\)](#), [DSC\(\)](#), [DSCClassifier\(\)](#), [DSOutlier\(\)](#), [DST\(\)](#), [DST\\_SlidingWindow\(\)](#), [DST\\_WriteStream\(\)](#), [evaluate](#), [predict\(\)](#), [stream\\_pipeline](#), [update\(\)](#)

Other DSRegressor: [DSRegressor\\_SlidingWindow\(\)](#)

**Examples**

```
DSClassifier()
```

---

```
DSRegressor_SlidingWindow
```

*DSRegressor\_SlidingWindow – Data Stream Regressor Using a Sliding Window*

---

**Description**

The Regressor keeps a sliding window for the stream and rebuilds a regression model at regular intervals. By default it builds a decision tree using [lm\(\)](#).

**Usage**

```
DSRegressor_SlidingWindow(formula, model = stats::lm, window, rebuild, ...)
```

**Arguments**

formula	a formula for the classification problem.
model	regression model (that has a formula interface).
window	size of the sliding window.
rebuild	interval (number of points) for rebuilding the regression. Set rebuild to Inf to prevent automatic rebuilding. Rebuilding can be initiated manually when calling <code>update()</code> .
...	additional parameters are passed on to the regressor (default is <code>lm()</code> ).

**Details**

This constructor creates a regressor based on [DST\\_SlidingWindow](#). The regressor has a `update()` and `predict()` method.

**Value**

An object of class `DST_SlidingWindow`.

**Author(s)**

Michael Hahsler

**See Also**

Other DSRegressor: [DSRegressor\(\)](#)

**Examples**

```
library(stream)

# create a data stream for the iris dataset
data <- iris[sample(nrow(iris)), ]
stream <- DSD_Memory(data)

# define the stream Regressor.
cl <- DSRegressor_SlidingWindow(
  Sepal.Length ~ Petal.Length + Petal.Length,
  window = 50,
  rebuild = 10
)
cl

# update the regressor with 100 points from the stream
update(cl, stream, 100)

# predict the class for the next 50 points
newdata <- get_points(stream, n = 50)
pr <- predict(cl, newdata)
pr

plot(pr, newdata$Sepal.Length)
abline(0, 1, col = "red")

# get the tree model
get_model(cl)
```

---

DST

*Conceptual Base Class for All Data Stream Mining Tasks*

---

**Description**

Conceptual base class for all data stream mining tasks.

**Usage**

```
DST(...)
```

```
description(x, ...)
```

```
get_model(x, ...)
```

**Arguments**

...	Further arguments.
x	an object of a concrete implementation of a DST.

## Details

Base class for data stream mining tasks. Types of DST are

- [DSAggregate](#) to aggregate data streams (e.g., with a sliding window).
- [DSC](#) for data stream clustering.
- [DSClassifier](#) classification for data streams.
- [DSRegressor](#) regression for data streams.
- [DSOutlier](#) outlier detection for data streams.
- [DSFP](#) frequent pattern mining for data streams.

The common interface for all [DST](#) classes consists of

- [update\(\)](#) update the DST with data points.
- [description\(\)](#) a string describing the DST.
- [get\\_model\(\)](#) returns the DST's current model (often as a `data.frame` or a R model object).
- [predict\(\)](#) use the learned DST model to make predictions.

and the methods in the Methods Section below.

## Author(s)

Michael Hahsler

## See Also

Other DST: [DSAggregate\(\)](#), [DSC\(\)](#), [DSClassifier\(\)](#), [DSOutlier\(\)](#), [DSRegressor\(\)](#), [DST\\_SlidingWindow\(\)](#), [DST\\_WriteStream\(\)](#), [evaluate](#), [predict\(\)](#), [stream\\_pipeline](#), [update\(\)](#)

## Examples

```
DST()
```

---

DST\_Multi

*Apply Multiple Task to the Same Data Stream*

---

## Description

Apply multiple task ([DST](#)) to the same data stream. The tasks can be accessed as a list as `$dsts`.

## Usage

```
DST_Multi(dsts)
```

## Arguments

`dsts` a list of [DST](#) objects.



**Author(s)**

Michael Hahsler

**Examples**

```

set.seed(1500)

stream <- DSD_Gaussians(k = 3, d = 2)

## define multiple tasks as a list
tasks <- DST_Multi(list(
  DSAggregate_Window(horizon = 10),
  DSC_DStream(gridsize = 0.1)
))
tasks

## update both tasks with the same stream
update(tasks, stream, n = 1000)

## inspect the results of the tasks
tasks$dsts[[1]]
get_model(tasks$dsts[[1]])

tasks$dsts[[2]]
plot(tasks$dsts[[2]])

```

---

DST\_SlidingWindow

*DST\_SlidingWindow – Call R Functions on a Sliding Window*


---

**Description**

Keeps a sliding window of the data stream and calls a function at regular intervals with the contents of the window.

**Usage**

```

DST_SlidingWindow(f, window, rebuild, ...)

## S3 method for class 'DST_SlidingWindow'
update(
  object,
  dsd,
  n = 1L,
  return = c("nothing", "model"),
  rebuild = FALSE,
  ...
)

## S3 method for class 'DST_SlidingWindow'
predict(object, newdata, ...)

```

**Arguments**

f	the function to be called.
window	size of the sliding window.
rebuild	logical; perform a rebuild after the update.
...	additional parameters passed on to the predict() function of the underlying model.
object	the updated DST_SlidingWindow object.
dsd	A <a href="#">DSD</a> object with the data stream.
n	number of points from dsd to use for the update.
return	a character string indicating what update returns. The default is "nothing" and "model" returns the aggregated data.
newdata	dataframe with the new data.

**Details**

Keeps a sliding window of the data stream and calls a function at regular intervals with the contents of the window. The function needs to have the form

```
f <- function(data, ...) { ... }
```

where data is the data.frame with the points in the sliding window (See [get\\_points\(\)](#) in [DSAggregate\\_Window](#)). The function will be executed at regular intervals after [update\(\)](#) was called with fixed number of points. [update\(..., rebuild = TRUE\)](#) can be used to force recomputing the function. This can be used with  $n = 0$  to recompute it even without adding more points. The last valid function result can be retrieved/

Many modelling functions provide a formula interface which lets them be directly used inside a DST\_SlidingWindow (see Examples section).

If the function returns a model that supports [predict\(\)](#), then [predict](#) can directly be called on the DST\_SlidingWindow object.

**Value**

An object of class DST\_SlidingWindow.

**Author(s)**

Michael Hahsler

**See Also**

Other DST: [DSAggregate\(\)](#), [DSC\(\)](#), [DSCClassifier\(\)](#), [DSOutlier\(\)](#), [DSRegressor\(\)](#), [DST\(\)](#), [DST\\_WriteStream\(\)](#), [evaluate](#), [predict\(\)](#), [stream\\_pipeline](#), [update\(\)](#)

**Examples**

```
library(stream)

# create a data stream for the iris dataset
data <- iris[sample(nrow(iris)), ]
stream <- DSD_Memory(data)
stream

## Example 1: Use a function on the sliding window
summarizer <- function(data) summary(data)

s <- DST_SlidingWindow(summarizer,
  window = 100, rebuild = 50)
s

# update window with 49 points. The function is not yet called
update(s, stream, 49)
get_model(s)

# updating with the 50th point will trigger a function call (see rebuild parameter)
# note that the window is only 1/2 full and we have 50 NAs
update(s, stream, 1)
get_model(s)

# 50 more points and the function will be recomputed
update(s, stream, 50)
get_model(s)

## Example 2: Use classifier on the sliding window
reset_stream(stream)

# rpart, like most models in R, already have a formula interface that uses a
# data parameter. We can use these types of models directly
library(rpart)
cl <- DST_SlidingWindow(
  rpart, formula = Species ~ Petal.Length + Petal.Width,
  window = 100, rebuild = 50)
cl

# update window with 50 points so the model is built
update(cl, stream, 50)
get_model(cl)

# update with 40 more points and force the function to be recomputed even though it would take
# 50 points to automatically rebuild.
update(cl, stream, 40, rebuild = TRUE)
get_model(cl)

# rpart supports predict, so we can use it directly with the DST_SlidingWindow
new_points <- get_points(stream, n = 5)
predict(cl, new_points, type = "class")
```

```
## Example 3: Regression using a sliding window
reset_stream(stream)

## lm can be directly used
reg <- DST_SlidingWindow(
  lm, formula = Sepal.Length ~ Petal.Width + Petal.Length,
  window = 100, rebuild = 50)
reg

update(reg, stream, 100)
get_model(reg)

# lm supports predict, so we can use it directly with the DST_SlidingWindow
new_points <- get_points(stream, n = 5)
predict(reg, new_points)
```

---

DST\_WriteStream

*Task to Write a Stream to a File or a Connection*


---

### Description

Writes points from a data stream DSD object to a file or a connection.

### Usage

```
DST_WriteStream(file, append = FALSE, ...)
```

```
## S3 method for class 'DST_WriteStream'
close_stream(dsd, ...)
```

### Arguments

file	A file name or a R connection to be written to.
append	Append the data to an existing file.
...	further arguments are passed on to <a href="#">write_stream()</a> . Note that close is always FALSE and cannot be specified.
dsd	a DST_WriteStream object with an open connection.

### Details

**Note:** header = TRUE is not supported for files. The header would be added for every call for update.

### Author(s)

Michael Hahsler

**See Also**

Other DST: [DSAggregate\(\)](#), [DSC\(\)](#), [DSCClassifier\(\)](#), [DSOutlier\(\)](#), [DSRegressor\(\)](#), [DST\(\)](#), [DST\\_SlidingWindow\(\)](#), [evaluate](#), [predict\(\)](#), [stream\\_pipeline](#), [update\(\)](#)

**Examples**

```
set.seed(1500)

stream <- DSD_Gaussians(k = 3, d = 2)
writer <- DST_WriteStream(file = "data.txt", info = TRUE)

update(writer, stream, n = 2)
readLines("data.txt")
update(writer, stream, n = 3)
readLines("data.txt")

# clean up
close_stream(writer)

file.remove("data.txt")
```

---

 evaluate

*Evaluate a Data Stream Mining Task*


---

**Description**

Generic function to calculate evaluation measures for a data stream mining task [DST](#) on a data stream [DSD](#) object.

**Usage**

```
evaluate_static(object, dsd, measure, n, ...)

evaluate_stream(object, dsd, measure, n, horizon, ..., verbose = FALSE)
```

**Arguments**

object	The <a href="#">DST</a> object that the evaluation measure is being requested from.
dsd	The <a href="#">DSD</a> object used to create the test data.
measure	Evaluation measure(s) to use. If missing then all available measures are returned.
n	The number of data points being requested.
...	Further arguments are passed on to the specific implementation (e.g., see <a href="#">evaluate.DSC</a> )
horizon	Evaluation is done using horizon many previous points (see detail section).
verbose	Report progress?

## Details

We define two generic evaluation functions:

- `evaluate_static()` evaluates the current [DST](#) model on new data without updating the model.
- `evaluate_stream()` evaluates the [DST](#) model using *prequential error estimation* (see Gama, Sebastiao and Rodrigues; 2013). The data points in the horizon are first used to calculate the evaluation measure and then they are used for updating the cluster model. A horizon of ‘ means that each point is evaluated and then used to update the model.

The available evaluation measures depend on the task. Currently available task to evaluate:

- [DSC](#) via `evaluate.DSC`

## Value

`evaluate` returns an object of class `stream_eval` which is a numeric vector of the values of the requested measures.

## Author(s)

Michael Hahsler

## References

Joao Gama, Raquel Sebastiao, Pedro Pereira Rodrigues (2013). On evaluating stream learning algorithms. *Machine Learning*, March 2013, Volume 90, Issue 3, pp 317-346.

## See Also

Other [DST](#): [DSAggregate\(\)](#), [DSC\(\)](#), [DSCClassifier\(\)](#), [DSOutlier\(\)](#), [DSRegressor\(\)](#), [DST\(\)](#), [DST\\_SlidingWindow\(\)](#), [DST\\_WriteStream\(\)](#), [predict\(\)](#), [stream\\_pipeline](#), [update\(\)](#)

Other evaluation: [animate\\_cluster\(\)](#), [evaluate.DSC](#)

---

evaluate.DSC

*Evaluate a Stream Clustering Task*

---

## Description

Methods for the generic functions `evaluate_static()` and `evaluate_stream()` to calculate evaluation measures for micro or macro-clusters created by a [DSC](#) on the a [DSD](#) object.

**Usage**

```
## S3 method for class 'DSC'
evaluate_static(
  object,
  dsd,
  measure,
  n = 100,
  type = c("auto", "micro", "macro"),
  assign = "micro",
  assignmentMethod = c("auto", "model", "nn"),
  excludeNoise = FALSE,
  callbacks = list(),
  ...
)
```

```
## S3 method for class 'DSC'
evaluate_stream(
  object,
  dsd,
  measure,
  n = 1000,
  horizon = 100,
  type = c("auto", "micro", "macro"),
  assign = "micro",
  assignmentMethod = c("auto", "model", "nn"),
  excludeNoise = FALSE,
  callbacks = NULL,
  ...,
  verbose = FALSE
)
```

**Arguments**

object	The <a href="#">DSC</a> object that the evaluation measure is being requested from.
dsd	The <a href="#">DSD</a> object that holds the initial training data for the DSC.
measure	Evaluation measure(s) to use. If missing then all available measures are returned.
n	The number of data points being requested.
type	Use micro- or macro-clusters for evaluation. Auto used the class of <a href="#">DSC</a> to decide.
assign	Assign points to micro or macro-clusters?
assignmentMethod	How are points assigned to clusters for evaluation (see <a href="#">predict()</a> )?
excludeNoise	logical; Should noise points in the data stream be excluded from the calculation?
callbacks	A named list of functions to calculate custom evaluation measures.
...	Unused arguments are ignored.

horizon	Evaluation is done using horizon many previous points (see detail section).
verbose	logical; Report progress?

## Details

For evaluation, each data point is assigned to its nearest cluster using Euclidean distance to the cluster centers. Then for each cluster the majority class is determined. Based on the majority class several evaluation measures can be computed.

We provide two evaluation methods:

- `evaluate_static()` evaluates the current static clustering using new data without updating the model.
- `evaluate_stream()` evaluates the clustering process using *prequential error estimation* (see Gama, Sebastiao and Rodrigues; 2013). The current model is first applied to the data points in the horizon to calculate the evaluation measures. Then, the cluster model is updated with the points.

## Evaluation Measures

Many evaluation measures are available using code from other packages including `cluster::silhouette()`, `clue::cl_agreement()`, and `fpc::cluster.stats()`.

The following information items are available:

- "numPoints" number of points used for evaluation.
- "numMicroClusters" number of micro-clusters
- "numMacroClusters" number of macro-clusters
- "numClasses" number of classes

The following noise-related/outlier items are available:

- "noisePredicted" Number data points predicted as noise
- "noiseActual" Number of data points which are actually noise
- "noisePrecision" Precision of the predicting noise (i.e., number of correctly predicted noise points over the total number of points predicted as noise)
- "outlierJaccard" - A variant of the Jaccard index used to assess outlier detection accuracy (see Krleza et al (2020)). Outlier Jaccard index is calculated as  $TP / (TP + FP + UNDETECTED)$ .

The following internal evaluation measures are available:

- "SSQ" within cluster sum of squares. Assigns each point to its nearest center from the clustering and calculates the sum of squares. Noise points in the data stream are always ignored.
- "silhouette" average silhouette width. Actual noise points which stay unassigned by the clustering algorithm are ignored; regular points that are unassigned by the clustering algorithm form their own noise cluster) (**cluster**)
- "average.between" average distance between clusters (**fpc**)
- "average.within" average distance within clusters (**fpc**)
- "max.diameter" maximum cluster diameter (**fpc**)



- "min.separation" minimum cluster separation (**fpc**)
- "ave.within.cluster.ss" a generalization of the within clusters sum of squares (half the sum of the within cluster squared dissimilarities divided by the cluster size) (**fpc**)
- "g2" Goodman and Kruskal's Gamma coefficient (**fpc**)
- "pearsongamma" correlation between distances and a 0-1-vector where 0 means same cluster, 1 means different clusters (**fpc**)
- "dunn" Dunn index (minimum separation / maximum diameter) (**fpc**)
- "dunn2" minimum average dissimilarity between two cluster / maximum average within cluster dissimilarity (**fpc**)
- "entropy" entropy of the distribution of cluster memberships (**fpc**)
- "wb.ratio" average.within/average.between (**fpc**)

The following external evaluation measures are available:

- "precision", "recall", "F1" F1. A true positive (TP) decision assigns two points in the same true cluster also to the same cluster, a true negative (TN) decision assigns two points from two different true clusters to two different clusters. A false positive (FP) decision assigns two points from the same true cluster to two different clusters. A false negative (FN) decision assigns two points from the same true cluster to different clusters.  

$$\text{precision} = TP / (TP + FP)$$

$$\text{recall} = TP / (TP + FN)$$

The F1 measure is the harmonic mean of precision and recall.
- "purity" Average purity of clusters. The purity of each cluster is the proportion of the points of the majority true group assigned to it (see Cao et al. (2006)).
- "classPurity" (of real clusters; see Wan et al (2009)).
- "fpr" false positive rate.
- "Euclidean" Euclidean dissimilarity of the memberships (see Dimitriadou, Weingessel and Hornik (2002)) (**clue**)
- "Manhattan" Manhattan dissimilarity of the memberships (**clue**)
- "Rand" Rand index (see Rand (1971)) (**clue**)
- "cRand" Adjusted Rand index (see Hubert and Arabie (1985)) (**clue**)
- "NMI" Normalized Mutual Information (see Strehl and Ghosh (2002)) (**clue**)
- "KP" Katz-Powell index (see Katz and Powell (1953)) (**clue**)
- "angle" maximal cosine of the angle between the agreements (**clue**) - "diag" maximal co-classification rate (**clue**)
- "FM" Fowlkes and Mallows's index (see Fowlkes and Mallows (1983)) (**clue**)
- "Jaccard" Jaccard index (**clue**)
- "PS" Prediction Strength (see Tibshirani and Walter (2005)) (**clue**) %
- "corrected.rand" corrected Rand index (**fpc**)
- "vi" variation of information (VI) index (**fpc**)

Many measures are the average over all clusters. For example, purity is the average purity over all clusters.

For `DSC_Micro` objects, data points are assigned to micro-clusters and then each micro-cluster is evaluated. For `DSC_Macro` objects, data points by default (assign = "micro") also assigned to micro-clusters, but these assignments are translated to macro-clusters. The evaluation is here done for macro-clusters. This is important when macro-clustering is done with algorithms which do not create spherical clusters (e.g. hierarchical clustering with single-linkage or DBSCAN) and this assignment to the macro-clusters directly (i.e., their center) does not make sense.

Using type and assign, the user can select how to assign data points and at what level (micro or macro) to evaluate.

`evaluate_cluster()` is used to evaluate an evolving data stream using the method described by Wan et al. (2009). Of the  $n$  data points horizon many points are clustered and then the evaluation measure is calculated on the same data points. The idea is to find out if the clustering algorithm was able to adapt to the changing stream.

### Custom Evaluation Measures

The parameter callbacks can be supplied with a named list with functions with the signature `function(actual, predict, points, centers, dsc)` as elements. See the Examples sections for details.

### Value

`evaluate` returns an object of class `stream_eval` which is a numeric vector of the values of the requested measures and two attributes, "type" and "assign", to see at what level the evaluation was done.

### Author(s)

Michael Hahsler, Matthew Bolanos, John Forrest, and Dalibor Krleža

### References

- Joao Gama, Raquel Sebastiao, Pedro Pereira Rodrigues (2013). On evaluating stream learning algorithms. *Machine Learning*, March 2013, Volume 90, Issue 3, pp 317-346.
- F. Cao, M. Ester, W. Qian, A. Zhou (2006). Density-Based Clustering over an Evolving Data Stream with Noise. *Proceeding of the 2006 SIAM Conference on Data Mining*, 326-337.
- E. Dimitriadou, A. Weingessel and K. Hornik (2002). A combination scheme for fuzzy clustering. *International Journal of Pattern Recognition and Artificial Intelligence*, 16, 901-912.
- E. B. Fowlkes and C. L. Mallows (1983). A method for comparing two hierarchical clusterings. *Journal of the American Statistical Association*, 78, 553-569.
- L. Hubert and P. Arabie (1985). Comparing partitions. *Journal of Classification*, 2, 193-218.
- W. M. Rand (1971). Objective criteria for the evaluation of clustering methods. *Journal of the American Statistical Association*, 66, 846-850.
- L. Katz and J. H. Powell (1953). A proposed index of the conformity of one sociometric measurement to another. *Psychometrika*, 18, 249-256.
- A. Strehl and J. Ghosh (2002). Cluster ensembles - A knowledge reuse framework for combining multiple partitions. *Journal of Machine Learning Research*, 3, 583-617.

R. Tibshirani and G. Walter (2005). Cluster validation by Prediction Strength. *Journal of Computational and Graphical Statistics*, 14/3, 511-528.

L. Wan, W.K. Ng, X.H. Dang, P.S. Yu and K. Zhang (2009). Density-Based Clustering of Data Streams at Multiple Resolutions, *ACM Transactions on Knowledge Discovery from Data*, 3(3).

D. Krleža, B. Vrdoljak, and M. Brčić (2020). Statistical Hierarchical Clustering Algorithm for Outlier Detection in Evolving Data Streams, *Springer Machine Learning*.

### See Also

`cluster::silhouette()`, `clue::cl_agreement()`, and `fpc::cluster.stats()`.

Other DSC: `DSC()`, `DSC_Macro()`, `DSC_Micro()`, `DSC_R()`, `DSC_SlidingWindow()`, `DSC_Static()`, `DSC_TwoStage()`, `animate_cluster()`, `get_assignment()`, `plot.DSC()`, `predict()`, `prune_clusters()`, `read_saveDSC`, `recluster()`

Other evaluation: `animate_cluster()`, `evaluate`

### Examples

```
# Example 1: Static Evaluation
set.seed(0)
stream <- DSD_Gaussians(k = 3, d = 2)

dstream <- DSC_DStream(gridsize = 0.05, Cm = 1.5)
update(dstream, stream, 500)
plot(dstream, stream)

# Evaluate the micro-clusters in the clustering
# Note: we use here only n = 100 points for evaluation to speed up execution
evaluate_static(dstream, stream, n = 100)

evaluate_static(dstream, stream,
  measure = c("numMicro", "numMacro", "purity", "crand", "SSQ"),
  n = 100)

# DStream also provides macro clusters. Evaluate macro clusters with type = "macro"
# Note that SSQ and cRand increase.
plot(dstream, stream, type = "macro")
evaluate_static(dstream, stream, type = "macro",
  measure = c("numMicro", "numMacro", "purity", "crand", "SSQ"),
  n = 100)

# Points are by default assigned to micro clusters using the method
# specified for the clustering algorithm.
# However, points can also be assigned to the closest macro-cluster using
# assign = "macro".
evaluate_static(dstream, stream, type = "macro", assign = "macro",
  measure = c("numMicro", "numMacro", "purity", "crand", "SSQ"),
  n = 100)

# Example 2: Evaluate with Noise/Outliers
stream <- DSD_Gaussians(k = 3, d = 2, noise = .05)
```

```

dstream <- DSC_DStream(gridsize = 0.05, Cm = 1.5)
update(dstream, stream, 500)

# For cRand, noise is its own group, for SSQ, actual noise is always
# excluded.
plot(dstream, stream, 500)
evaluate_static(dstream, stream, n = 100,
  measure = c("numPoints", "noisePredicted", "noiseActual",
    "noisePrecision", "outlierJaccard", "cRand", "SSQ"))

# Note that if noise is excluded, the number of used points is reduced.
evaluate_static(dstream, stream, n = 100,
  measure = c("numPoints", "noisePredicted", "noiseActual",
    "noisePrecision", "outlierJaccard", "cRand", "SSQ"), excludeNoise = TRUE)

# Example 3: Evaluate an evolving data stream
stream <- DSD_Benchmark(1)
dstream <- DSC_DStream(gridsize = 0.05, lambda = 0.1)

evaluate_stream(dstream, stream, type = "macro", assign = "micro",
  measure = c("numMicro", "numMacro", "purity", "cRand"),
  n = 600, horizon = 100)

if (interactive()){
# animate the clustering process
reset_stream(stream)
dstream <- DSC_DStream(gridsize = 0.05, lambda = 0.1)
animate_cluster(dstream, stream, horizon = 100, n = 5000,
  measure = "cRand", type = "macro", assign = "micro",
  plot.args = list(type = "both", xlim = c(0,1), ylim = c(0,1)))
}

# Example 4: Add a custom measure as a callback
callbacks <- list(
  noisePercentage = function(actual, predict, points, centers, dsc) {
    sum(actual == 0L) / length(actual)
  },
  noiseFN = function(actual, predict, points, centers, dsc) {
    sum(actual == 0L & predict != 0L)
  },
  noiseFP = function(actual, predict, points, centers, dsc) {
    sum(actual != 0L & predict == 0L)
  }
)

stream <- DSD_Gaussians(k = 3, d = 2, noise = .2)
dstream <- DSC_DStream(gridsize = 0.05, Cm = 1.5)
update(dstream, stream, 500)

evaluate_static(dstream, stream,
  measure = c("numPoints", "noiseActual", "noisePredicted",
    "noisePercentage", "noiseFN", "noiseFP"),

```

```

    callbacks = callbacks, n = 100)

evaluate_static(dstream, stream, callbacks = callbacks)

```

---

get_assignment	Assignment Data Points to Clusters <i>deprecated</i>
----------------	--

---

## Description

**Deprecation Notice:** use `predict()` for a more general interface to apply a data stream model to new data. `get_assignment()` is deprecated.

## Usage

```

get_assignment(
  dsc,
  points,
  type = c("auto", "micro", "macro"),
  method = "auto",
  ...
)

## S3 method for class 'DSC'
get_assignment(
  dsc,
  points,
  type = c("auto", "micro", "macro"),
  method = c("auto", "nn", "model"),
  ...
)

```

## Arguments

dsc	The <a href="#">DSC</a> object with the clusters for assignment.
points	The points to be assigned as a data.frame.
type	Use micro- or macro-clusters in <a href="#">DSC</a> for assignment.
method	assignment method <ul style="list-style-type: none"> <li>• "model" uses the assignment method of the underlying algorithm (unassigned points return NA). Not all algorithms implement this option.</li> <li>• "nn" performs nearest neighbor assignment using Euclidean distance.</li> <li>• "auto" uses the model assignment method. If this method is not implemented/available then method "nn" is used instead.</li> </ul>
...	Additional arguments are passed on.

**Details**

Get the assignment of data points to clusters in a DSC using the model's assignment rules or nearest neighbor assignment. The clustering is not modified.

Each data point is assigned either using the original model's assignment rule or Euclidean nearest neighbor assignment. If the user specifies the model's assignment strategy, but is not available, then nearest neighbor assignment is used and a warning is produced.

**Value**

A vector containing the assignment of each point. NA means that a data point was not assigned to a cluster.

**Author(s)**

Michael Hahsler

**See Also**

Other DSC: [DSC\(\)](#), [DSC\\_Macro\(\)](#), [DSC\\_Micro\(\)](#), [DSC\\_R\(\)](#), [DSC\\_SlidingWindow\(\)](#), [DSC\\_Static\(\)](#), [DSC\\_TwoStage\(\)](#), [animate\\_cluster\(\)](#), [evaluate.DSC](#), [plot.DSC\(\)](#), [predict\(\)](#), [prune\\_clusters\(\)](#), [read\\_saveDSC](#), [recluster\(\)](#)

**Examples**

```
stream <- DSD_Gaussians(k = 3, d = 2, noise = .05)

dbstream <- DSC_DBSTREAM(r = .1)
update(dbstream, stream, n = 100)

# find the assignment for the next 100 points to
# micro-clusters in dsc. This uses the model's assignment function
points <- get_points(stream, n = 100)
a <- predict(dbstream, points)
head(a)

# show the MC assignment areas. Assigned points as blue circles and
# the unassigned points as red dots
plot(dbstream, stream, assignment = TRUE, type = "none")
points(points[!is.na(a[, ".class"]),], col = "blue")
points(points[is.na(a[, ".class"]),], col = "red", pch = 20)

# use nearest neighbor assignment instead
a <- predict(dbstream, points, method = "nn")
head(a)
```

**Description**

Gets points from a [DSD](#) object.

**Usage**

```
get_points(x, ...)  
  
## S3 method for class 'DSD'  
get_points(x, n = 1L, info = TRUE, ...)  
  
remove_info(points)
```

**Arguments**

x	A <a href="#">DSD</a> object.
...	Additional parameters to pass to the <code>get_points()</code> implementations.
n	integer; request up to n points from the stream. n = -1 returns all remaining points from limited streams.
info	return additional columns with information about the data point (e.g., a known cluster assignment).
points	a data.frame with points.

**Details**

Each DSD object has a unique way for creating/returning data points, but they all are called through the generic function, `get_points()`. This is done by using the S3 class system. See the man page for the specific [DSD](#) class on the semantics for each implementation of `get_points()`.

**Additional Point Information**

Additional point information (e.g., known cluster/class assignment, noise status) can be requested with `info = TRUE`. This information is returned as additional columns. The column names start with `.` and are ignored by [DST](#) implementations. `remove_info()` is a convenience function to remove the information columns. Examples are

- `.id` for point IDs
- `.class` for known cluster/class labels used for plotting and evaluation
- `.time` a time stamp for the point (can be in seconds or an index for ordering)

**Resetting a Stream**

Many streams can be reset using `reset_stream()`.

**Value**

Returns a [data.frame](#) with (up to) n rows and as many columns as x produces.

**Author(s)**

Michael Hahsler

**See Also**

Other DSD: [DSD\(\)](#), [DSD\\_BarsAndGaussians\(\)](#), [DSD\\_Benchmark\(\)](#), [DSD\\_Cubes\(\)](#), [DSD\\_Gaussians\(\)](#), [DSD\\_MG\(\)](#), [DSD\\_Memory\(\)](#), [DSD\\_Mixture\(\)](#), [DSD\\_NULL\(\)](#), [DSD\\_ReadDB\(\)](#), [DSD\\_ReadStream\(\)](#), [DSD\\_Target\(\)](#), [DSD\\_UniformNoise\(\)](#), [DSD\\_mlbenchData\(\)](#), [DSD\\_mlbenchGenerator\(\)](#), [DSF\(\)](#), [animate\\_data\(\)](#), [close\\_stream\(\)](#), [plot.DSD\(\)](#), [reset\\_stream\(\)](#)

**Examples**

```
stream <- DSD_Gaussians()
points <- get_points(stream, n = 5)
points

remove_info(points)
```

---

MGC

*Moving Generator Cluster*

---

**Description**

Creates an evolving cluster for use as a component of a [DSD\\_MG](#) data stream.

**Usage**

```
MGC(...)
```

```
MGC_Function(density, center, parameter, shape = Shape_Gaussian)
```

```
MGC_Linear(dimension = 2, keyframelist = NULL, shape = Shape_Gaussian)
```

```
keyframe(time, density, center, parameter, reset = FALSE)
```

```
add_keyframe(x, time, density, center, parameter, reset = FALSE)
```

```
get_keyframes(x)
```

```
remove_keyframe(x, time)
```

```
MGC_Noise(density, range)
```



```
MGC_Random(density, center, parameter, randomness = 1, shape = Shape_Gaussian)
```

```
Shape_Gaussian(center, parameter)
```

```
Shape_Block(center, parameter)
```

```
MGC_Static(density = 1, center, parameter, shape = Shape_Gaussian)
```

### Arguments

...	Further arguments.
density	The density of the cluster. For 'MGC_Function, this attribute is a function and defines the density of a cluster (i.e., how many points it creates) at each given timestamp.
center	A list that defines the center of the cluster. The list should have a length equal to the dimensionality. For MGC_Function, this list consists of functions that define the movement of the cluster. For MGC_Random, this attribute defines the beginning location for the MGC before it begins moving.
parameter	Parameters for the shape. For the default shape Shape_Gaussian the parameter is the standard deviation, one per dimension. If a single value is specified then it is recycled for all dimensions.
shape	A function creating the shape of the cluster. It gets passed on the parameters argument from above. Available functions are Shape_Gaussian (the parameters are a vector containing standard deviations) and Shape_Block (parameters are the dimensions of the uniform block).
dimension	Dimensionality of the data stream.
keyframelist	a list of keyframes to initialize the MGC_Linear object with.
time	The time stamp the keyframe should be located or which keyframe should be removed.
reset	Should the cluster reset to the first keyframe (time 0) after this keyframe is finished?
x	An object of class MGC_Linear.
range	The area in which the noise should appear.
randomness	The maximum amount the cluster will move during one time step.

### Details

An MGC describes a single cluster for use as a component in a [DSD\\_MG](#). Different [MGCs](#) allow the user to express different cluster behaviors within a single data stream. Static, (i.e., not moving) clusters are defined as:

- MGC\_Static cluster positions are fixed
- MGC\_Noise allows to add random noise.

Moving (evolving) clusters are defined as:

- MGC\_Linear creates an evolving cluster for a who's behavior is determined by keyframes. Several keyframe functions are provided to create, add and remove keyframes. See Examples section for details.
- MGC\_Function allows to specify density, center, and parameter as a function of time.
- MGC\_Random allows for a creation of a cluster that follows a random walk.

Cluster shapes can be specified using the functions:

- Shape\_Gaussian
- Shape\_Block

New shapes can be defined as a function with parameters center and parameter that return a single new point. Here is an example:

```
Shape_Gaussian <- function(center, parameter)
  rnorm(length(center), mean = center, sd = parameter)
```

### Author(s)

Matthew Bolanos

### See Also

[DSD\\_MG](#) for details on how to use an MGC within a [DSD](#).

### Examples

```
MGC()

### Two static clusters (Gaussian with sd of .1 and a Block with width .4)
### with added noise
stream <- DSD_MG(dim = 2,
  MGC_Static(den = .45, center = c(1, 0), par = .1, shape = Shape_Gaussian),
  MGC_Static(den = .45, center = c(2, 0), par = .4, shape = Shape_Block),
  MGC_Noise( den = .1, range = rbind(c(0, 3), c(-1,1)))
)
stream

plot(stream)

### Example of several MGC_Randoms which define clusters that randomly move.
stream <- DSD_MG(dim = 2,
  MGC_Random(den = 100, center=c(1, 0), par = .1, rand = .2),
  MGC_Random(den = 100, center=c(2, 0), par = .4, shape = Shape_Block, rand = .2)
)

## Not run:
  animate_data(stream, 2500, xlim = c(0,3), ylim = c(-1,1), horizon = 100)

## End(Not run)
```

```

### Example of several MGC_Functions

### a block-shaped cluster moving from bottom-left to top-right increasing size
c1 <- MGC_Function(
  density = function(t){ 100 },
  parameter = function(t){ 1 * t },
  center = function(t) c(t, t),
  shape = Shape_Block
)

### a cluster moving in a circle (default shape is Gaussian)
c2 <- MGC_Function(
  density = function(t){ 25 },
  parameter = function(t){ 5 },
  center= function(t) c(sin(t / 10) * 50 + 50, cos(t / 10) * 50 + 50)
)

stream <- DSD_MG(dim = 2, c1, c2)

## adding noise after the stream was created
add_cluster(stream, MGC_Noise(den = 10, range = rbind(c(-20, 120), c(-20, 120))))

stream

## Not run:
animate_data(stream, 10000, xlim = c(-20, 120), ylim = c(-20, 120), horizon = 100)

## End(Not run)

### Example of several MGC_Linear: A single cluster splits at time 50 into two.
### Note that c2 starts at time = 50!
stream <- DSD_MG(dim = 2)
c1 <- MGC_Linear(dim = 2)
add_keyframe(c1, time = 1, dens = 50, par = 5, center = c(0, 0))
add_keyframe(c1, time = 50, dens = 50, par = 5, center = c(50, 50))
add_keyframe(c1, time = 100, dens = 50, par = 5, center = c(50, 100))
add_cluster(stream, c1)

c2 <- MGC_Linear(dim = 2, shape = Shape_Block)
add_keyframe(c2, time = 50, dens = 25, par = c(10, 10), center = c(50, 50))
add_keyframe(c2, time = 100, dens = 25, par = c(30, 30), center = c(100, 50))
add_cluster(stream, c2)

## Not run:
animate_data(stream, 5000, xlim = c(0, 100), ylim = c(0, 100), horiz = 100)

## End(Not run)

### two fixed and a moving cluster
stream <- DSD_MG(dim = 2,
  MGC_Static(dens = 1, par = .1, center = c(0, 0)),
  MGC_Static(dens = 1, par = .1, center = c(1, 1)),

```

```

MGC_Linear(dim = 2, list(
  keyframe(time = 0, dens = 1, par = .1, center = c(0, 0)),
  keyframe(time = 1000, dens = 1, par = .1, center = c(1, 1)),
  keyframe(time = 2000, dens = 1, par = .1, center = c(0, 0), reset = TRUE)
)))

noise <- MGC_Noise(dens = .1, range = rbind(c(-.2, 1.2), c(-.2, 1.2)))
add_cluster(stream, noise)

## Not run:
animate_data(stream, n = 2000 * 3.1, xlim = c(-.2, 1.2), ylim = c(-.2, 1.2), horiz = 200)

## End(Not run)

```

---

plot.DSC

*Plot Results of a Data Stream Clustering*


---

## Description

Method to plot the result of data stream data clustering. To plot [DSD](#) see [plot.DSD\(\)](#).

## Usage

```

## S3 method for class 'DSC'
plot(
  x,
  dsd = NULL,
  n = 500,
  col_points = NULL,
  col_clusters = c("red", "blue", "green"),
  weights = TRUE,
  scale = c(1, 5),
  cex = 1,
  pch = NULL,
  method = c("pairs", "scatter", "pca"),
  dim = NULL,
  type = c("auto", "micro", "macro", "both", "none"),
  assignment = FALSE,
  transform = NULL,
  ...
)

```

## Arguments

x	the <a href="#">DSC</a> object to be plotted.
dsd	a <a href="#">DSD</a> object to plot the data in the background.
n	number of plots taken from dsd to plot.

col_points, col_clusters	colors used for plotting.
weights	if TRUE then the cluster weight is used for symbol size. Alternatively, a vector with the size of the symbols for micro- and macro-clusters can be supplied.
scale	range for the symbol sizes used.
cex	size factor for symbols.
pch	symbol type for points.
method	method used for plotting: "pairs" (pairs plot), "scatter" (scatter plot), "pca" (plot first 2 principal components).
dim	an integer vector with the dimensions to plot. If NULL then for methods pairs and "pca" all dimensions are used and for "scatter" the first two dimensions are plotted.
type	Plot micro clusters (type = "micro"), macro clusters (type = "macro"), both micro and macro clusters (type = "both").
assignment	logical; show assignment area of micro-clusters.
transform	a function that maps data stream points onto a 2-D plane for plotting.
...	further arguments are passed on to <code>graphics::plot.default()</code> or <code>graphics::pairs()</code> . <b>graphics.</b>

**Author(s)**

Michael Hahsler

**See Also**

Other DSC: `DSC()`, `DSC_Macro()`, `DSC_Micro()`, `DSC_R()`, `DSC_SlidingWindow()`, `DSC_Static()`, `DSC_TwoStage()`, `animate_cluster()`, `evaluate.DSC`, `get_assignment()`, `predict()`, `prune_clusters()`, `read_saveDSC`, `recluster()`

Other plot: `animate_cluster()`, `animate_data()`, `plot.DSD()`

**Examples**

```
stream <- DSD_Gaussians(k = 3, d = 3, noise = 0.05)
```

```
## Example 1: Plot data
plot(stream)
```

```
## Example 2: Plot a clustering
dstream <- DSC_DStream(gridsize = 0.1)
update(dstream, stream, 500)
dstream
plot(dstream, stream)
```

```
## plot micro or macro-clusters only
plot(dstream, stream, type = "micro")
plot(dstream, stream, type = "macro")
```

```

## plot projected on the first two principal components
## and on dimensions 2 and 3
plot(dstream, stream, method = "pca")
plot(dstream, stream, dim = c(2, 3))

## D-Stream has a special implementation to show assignment areas
plot(dstream, stream, assignment = TRUE)

## Example 4: Use a custom transformation for plotting.
##   We fit PCA using 100 points and create a transformation
##   function to project the stream to the first two PCs.
pr <- princomp(get_points(stream, n = 100, info = FALSE))
trans <- function(x) predict(pr, x)[, 1:2 , drop = FALSE]

trans(get_points(stream, n = 3))

plot(dstream, stream, transform = trans)

```

---

plot.DSD

*Plot Data Stream Data*


---

## Description

Method to plot data stream data. To plot [DSC](#) see [plot.DSC\(\)](#).

## Usage

```

## S3 method for class 'DSD'
plot(
  x,
  n = 500,
  col = NULL,
  pch = NULL,
  ...,
  method = c("pairs", "scatter", "pca", "ts"),
  dim = NULL,
  alpha = 0.6,
  transform = NULL
)

```

## Arguments

x	the <a href="#">DSD</a> object to be plotted.
n	number of plots taken from x to plot.
col	colors used for points.
pch	symbol type.
...	further arguments are passed on to <a href="#">graphics::plot.default()</a> or <a href="#">graphics::pairs()</a> .

method	method used for plotting: "pairs" (pairs plot), "scatter" (scatter plot), "pca" (plot first 2 principal components), or "ts" (time series).
dim	an integer vector with the dimensions to plot. If NULL then for methods pairs and "pca" all dimensions are used and for "scatter" the first two dimensions are plotted.
alpha	alpha shading used to plot the points.
transform	a function that maps data stream points onto a 2-D plane for plotting.

**Author(s)**

Michael Hahsler

**See Also**

Other DSD: [DSD\(\)](#), [DSD\\_BarsAndGaussians\(\)](#), [DSD\\_Benchmark\(\)](#), [DSD\\_Cubes\(\)](#), [DSD\\_Gaussians\(\)](#), [DSD\\_MG\(\)](#), [DSD\\_Memory\(\)](#), [DSD\\_Mixture\(\)](#), [DSD\\_NULL\(\)](#), [DSD\\_ReadDB\(\)](#), [DSD\\_ReadStream\(\)](#), [DSD\\_Target\(\)](#), [DSD\\_UniformNoise\(\)](#), [DSD\\_mlbenchData\(\)](#), [DSD\\_mlbenchGenerator\(\)](#), [DSF\(\)](#), [animate\\_data\(\)](#), [close\\_stream\(\)](#), [get\\_points\(\)](#), [reset\\_stream\(\)](#)

Other plot: [animate\\_cluster\(\)](#), [animate\\_data\(\)](#), [plot.DSC\(\)](#)

**Examples**

```
stream <- DSD_Gaussians(k=3, d=3)

## plot data
plot(stream, n = 500)
plot(stream, method = "pca", n = 500)
plot(stream, method = "scatter", dim = c(1, 3), n = 500)

## create and plot micro-clusters
dstream <- DSC_DStream(gridsize = 0.1)
update(dstream, stream, 500)
plot(dstream)

## plot with data, projected on the first two principal components
## and dimensions 2 and 3
plot(dstream, stream)
plot(dstream, stream, method = "pca")
plot(dstream, stream, dim = c(2, 3))

## plot micro and macro-clusters
plot(dstream, stream, type = "both")

## plot a time series using the AirPassenger data with the total monthly
## passengers from 1949 to 1960) as a stream
AirPassengers
stream <- DSD_Memory(data.frame(
  .time = time(AirPassengers),
  passengers = AirPassengers))
```

```
get_points(stream, n = 10)
plot(stream, n = 100, method = "ts")
```

---

predict

*Make a Prediction for a Data Stream Mining Task*

---

## Description

predict() for data stream mining tasks [DST](#).

## Usage

```
## S3 method for class 'DST'
predict(object, newdata, ...)

## S3 method for class 'DSC'
predict(
  object,
  newdata,
  type = c("auto", "micro", "macro"),
  method = "auto",
  ...
)
```

## Arguments

object	The <a href="#">DST</a> object.
newdata	The points to make predictions for as a data.frame.
...	Additional arguments are passed on.
type	Use micro- or macro-clusters in <a href="#">DSC</a> for assignment.
method	assignment method <ul style="list-style-type: none"> <li>• "model" uses the assignment method of the underlying algorithm (unassigned points return NA). Not all algorithms implement this option.</li> <li>• "nn" performs nearest neighbor assignment using Euclidean distance.</li> <li>• "auto" uses the model assignment method. If this method is not implemented/available then method "nn" is used instead.</li> </ul>

## Value

A data.frame with columns containing the predictions. The columns depend on the type of the data stream mining task.

## Author(s)

Michael Hahsler



**See Also**

Other DST: [DSAggregate\(\)](#), [DSC\(\)](#), [DSCClassifier\(\)](#), [DSOutlier\(\)](#), [DSRegressor\(\)](#), [DST\(\)](#), [DST\\_SlidingWindow\(\)](#), [DST\\_WriteStream\(\)](#), [evaluate](#), [stream\\_pipeline](#), [update\(\)](#)

Other DSC: [DSC\(\)](#), [DSC\\_Macro\(\)](#), [DSC\\_Micro\(\)](#), [DSC\\_R\(\)](#), [DSC\\_SlidingWindow\(\)](#), [DSC\\_Static\(\)](#), [DSC\\_TwoStage\(\)](#), [animate\\_cluster\(\)](#), [evaluate.DSC](#), [get\\_assignment\(\)](#), [plot.DSC\(\)](#), [prune\\_clusters\(\)](#), [read\\_saveDSC](#), [recluster\(\)](#)

**Examples**

```
set.seed(1500)
stream <- DSD_Gaussians(k = 3, d = 2, noise = .1)

dbstream <- DSC_DBSTREAM(r = .1)
update(dbstream, stream, n = 100)
plot(dbstream, stream, type = "both")

# find the assignment for the next 100 points to
# micro-clusters in dsc. This uses the model's assignment function
points <- get_points(stream, n = 10)
points

pr <- predict(dbstream, points, type = "macro")
pr

# Note that the clusters are labeled in arbitrary order. Check the
# agreement.
agreement(pr[, ".class"], points[, ".class"])
```

---

prune\_clusters

*Prune Clusters from a Clustering*


---

**Description**

Creates a (static) copy of a clustering where a fraction of the weight or the number of clusters with the lowest weights were pruned.

**Usage**

```
prune_clusters(dsc, threshold = 0.05, weight = TRUE)
```

**Arguments**

dsc	The DSC object to be pruned.
threshold	The numeric vector of probabilities for the quantile.
weight	should a fraction of the total weight in the clustering be pruned? Otherwise a fraction of clusters is pruned.

**Value**

Returns an object of class DSC\_Static.

**Author(s)**

Michael Hahsler

**See Also**

[DSC\\_Static](#)

Other DSC: [DSC\(\)](#), [DSC\\_Macro\(\)](#), [DSC\\_Micro\(\)](#), [DSC\\_R\(\)](#), [DSC\\_SlidingWindow\(\)](#), [DSC\\_Static\(\)](#), [DSC\\_TwoStage\(\)](#), [animate\\_cluster\(\)](#), [evaluate.DSC](#), [get\\_assignment\(\)](#), [plot.DSC\(\)](#), [predict\(\)](#), [read\\_saveDSC](#), [recluster\(\)](#)

**Examples**

```
# 3 clusters with 10% noise
stream <- DSD_Gaussians(k=3, noise=0.1)

dbstream <- DSC_DBSTREAM(r=0.1)
update(dbstream, stream, 500)
dbstream
plot(dbstream, stream)

# prune lightest micro-clusters for 20% of the weight of the clustering
static <- prune_clusters(dbstream, threshold=0.2)
static
plot(static, stream)
```

---

read\_saveDSC

*Save and Read DSC Objects*

---

**Description**

Save and Read DSC objects safely (serializes the underlying data structure). This also works for **streamMOA** DSC objects.

**Usage**

```
saveDSC(object, file, ...)
```

```
readDSC(file)
```

**Arguments**

object	a DSC object.
file	filename.
...	further arguments.

**Author(s)**

Michael Hahsler

**See Also**[saveRDS](#) and [readRDS](#).

Other DSC: [DSC\(\)](#), [DSC\\_Macro\(\)](#), [DSC\\_Micro\(\)](#), [DSC\\_R\(\)](#), [DSC\\_SlidingWindow\(\)](#), [DSC\\_Static\(\)](#), [DSC\\_TwoStage\(\)](#), [animate\\_cluster\(\)](#), [evaluate.DSC](#), [get\\_assignment\(\)](#), [plot.DSC\(\)](#), [predict\(\)](#), [prune\\_clusters\(\)](#), [recluster\(\)](#)

**Examples**

```
stream <- DSD_Gaussians(k = 3, noise = 0.05)

# create clusterer with r = 0.05
dbstream1 <- DSC_DBSTREAM(r = .05)
update(dbstream1, stream, 1000)
dbstream1

saveDSC(dbstream1, file="dbstream.Rds")

dbstream2 <- readDSC("dbstream.Rds")
dbstream2

## cleanup
unlink("dbstream.Rds")
```

recluster

*Re-clustering micro-clusters***Description**

Use an **\*offline** macro clustering algorithm to recluster micro-clusters into a final clusters.

**Usage**

```
recluster(macro, micro, type = "auto", ...)
```

## S3 method for class 'DSC\_Macro'

```
recluster(macro, micro, type = "auto", ...)
```

**Arguments**

macro	an empty <a href="#">DSC_Macro</a> .
micro	an updated <a href="#">DSC_Micro</a> with micro-clusters.
type	controls which clustering is used from micro. Typically auto.
...	additional arguments passed on.

**Details**

Takes centers and weights of the micro-clusters and applies the macro clustering algorithm.

See [DSC\\_TwoStage](#) for a convenient combination of micro and macro clustering.

**Value**

The object macro is altered in place and contains the clustering.

**Author(s)**

Michael Hahsler

**See Also**

Other DSC: [DSC\(\)](#), [DSC\\_Macro\(\)](#), [DSC\\_Micro\(\)](#), [DSC\\_R\(\)](#), [DSC\\_SlidingWindow\(\)](#), [DSC\\_Static\(\)](#), [DSC\\_TwoStage\(\)](#), [animate\\_cluster\(\)](#), [evaluate.DSC](#), [get\\_assignment\(\)](#), [plot.DSC\(\)](#), [predict\(\)](#), [prune\\_clusters\(\)](#), [read\\_saveDSC](#)

**Examples**

```
set.seed(0)
### create a data stream and a micro-clustering
stream <- DSD_Gaussians(k = 3, d = 3)

### sample can be seen as a simple online clusterer where the sample points
### are the micro clusters.
sample <- DSC_Sample(k = 50)
update(sample, stream, 500)
sample

### recluster using k-means
kmeans <- DSC_Kmeans(k = 3)
recluster(kmeans, sample)

### plot clustering
plot(kmeans, stream, type = "both", main = "Macro-clusters (Sampling + k-means)")
```

---

reset\_stream

*Reset a Data Stream to its Beginning*

---

**Description**

Resets the position in a [DSD](#) object to the beginning or, if available, any other position in the stream.

**Usage**

```
reset_stream(dsd, pos = 1)
```

**Arguments**

dsd	An object of class a subclass of <a href="#">DSD</a> which implements a reset function.
pos	Position in the stream (the beginning of the stream is position 1).

**Details**

Resets the counter of the stream object. For example, for [DSD\\_Memory](#), the counter stored in the environment variable is moved back to 1. For [DSD\\_ReadCSV](#) objects, this is done by calling [seek\(\)](#) on the underlying connection.

reset\_stream() is implemented for:

- [DSD](#)
- [DSD\\_MG](#)
- [DSD\\_Memory](#)
- [DSD\\_ReadStream](#)
- [DSF](#)
- [DSF\\_Convolve](#)

**Author(s)**

Michael Hahsler

**See Also**

Other DSD: [DSD\(\)](#), [DSD\\_BarsAndGaussians\(\)](#), [DSD\\_Benchmark\(\)](#), [DSD\\_Cubes\(\)](#), [DSD\\_Gaussians\(\)](#), [DSD\\_MG\(\)](#), [DSD\\_Memory\(\)](#), [DSD\\_Mixture\(\)](#), [DSD\\_NULL\(\)](#), [DSD\\_ReadDB\(\)](#), [DSD\\_ReadStream\(\)](#), [DSD\\_Target\(\)](#), [DSD\\_UniformNoise\(\)](#), [DSD\\_m1benchData\(\)](#), [DSD\\_m1benchGenerator\(\)](#), [DSF\(\)](#), [animate\\_data\(\)](#), [close\\_stream\(\)](#), [get\\_points\(\)](#), [plot.DSD\(\)](#)

**Examples**

```
# initializing the objects
stream <- DSD_Gaussians()
replayer <- DSD_Memory(stream, 100)
replayer

p <- get_points(replayer, 50)
replayer

# reset replayer to the beginning of the stream
reset_stream(replayer)
replayer

# set replayer to position 21
reset_stream(replayer, pos = 21)
replayer
```

---

stream\_pipeline      *Create a Data Stream Pipeline*

---

## Description

Define a complete data stream pipe line consisting of a data stream, filters and a data mining task using %>%.

## Usage

```
DST_Runner(dsd, dst)
```

## Arguments

`dsd`                    A data stream (subclass of [DSD](#)) typically provided using a %>% (pipe).  
`dst`                    A data stream mining task (subclass of [DST](#)).

## Details

A data stream pipe line consisting of a data stream, filters and a data mining task:

```
DSD %>% DSF %>% DST_Runner
```

Once the pipeline is defined, it can be run using [update\(\)](#) where points are taken from the [DSD](#) data stream source, filtered through a sequence of [DSF](#) filters and then used to update the [DST](#) task.

[DST\\_Multi](#) can be used to update multiple models in the pipeline with the same stream.

## Author(s)

Michael Hahsler

## See Also

Other DST: [DSAggregate\(\)](#), [DSC\(\)](#), [DSCClassifier\(\)](#), [DSOutlier\(\)](#), [DSRegressor\(\)](#), [DST\(\)](#), [DST\\_SlidingWindow\(\)](#), [DST\\_WriteStream\(\)](#), [evaluate](#), [predict\(\)](#), [update\(\)](#)

## Examples

```
set.seed(1500)

# Set up a pipeline with a DSD data source, DSF Filters and then a DST task
cluster_pipeline <- DSD_Gaussians(k = 3, d = 2) %>%
  DSF_Scale() %>%
  DST_Runner(DSC_DBSTREAM(r = .3))

cluster_pipeline

# the DSD and DST can be accessed directly
cluster_pipeline$dsd
```

```

cluster_pipeline$dst

# update the DST using the pipeline, by default update returns the micro clusters
update(cluster_pipeline, n = 1000)

cluster_pipeline$dst
get_centers(cluster_pipeline$dst, type = "macro")
plot(cluster_pipeline$dst)

```

---

update	<i>Update a Data Stream Mining Task Model with Points from a Stream</i>
--------	---

---

### Description

update() for data stream mining tasks [DST](#).

### Usage

```

## S3 method for class 'DST'
update(object, dsd, n = 1L, return = "nothing", ...)

```

### Arguments

object	The <a href="#">DST</a> object.
dsd	A <a href="#">DSD</a> object with the data stream.
n	number of points from dsd to use for the update. Some DSD dsd accept n = -1 to update with all remaining points in the stream.
return	a character string indicating what update returns. The default is "nothing". Other possible values depend on the DST. Examples are "data", "model" and "assignment".
...	Additional arguments are passed on.

### Value

NULL or a data.frame n rows containing update information for each data point.

### Author(s)

Michael Hahsler

### See Also

Other DST: [DSAggregate\(\)](#), [DSC\(\)](#), [DSCClassifier\(\)](#), [DSOutlier\(\)](#), [DSRegressor\(\)](#), [DST\(\)](#), [DST\\_SlidingWindow\(\)](#), [DST\\_WriteStream\(\)](#), [evaluate](#), [predict\(\)](#), [stream\\_pipeline](#)

**Examples**

```

set.seed(1500)
stream <- DSD_Gaussians(k = 3, d = 2, noise = .1)

dbstream <- DSC_DBSTREAM(r = .1)
assignment <- update(dbstream, stream, n = 100, return = "assignment")
plot(dbstream, stream, type = "both")

# DBSTREAM returns cluster assignments (see DSC_DBSTREAM).
head(assignment)

```

---

write\_stream

*Write a Data Stream to a File*


---

**Description**

Writes points from a data stream DSD object to a file or a connection.

**Usage**

```

write_stream(
  dsd,
  file,
  n,
  block = 100000L,
  info = FALSE,
  append = FALSE,
  sep = ",",
  header = FALSE,
  row.names = FALSE,
  close = TRUE,
  ...
)

```

**Arguments**

dsd	The DSD object that will generate the data points for output.
file	A file name or a R connection to be written to.
n	The number of data points to be written. For finite streams, n = -1 writes all available data points.
block	Write stream in blocks to improve file I/O speed.
info	Save the class/cluster labels and other information columns with the data.
append	Append the data to an existing file. If FALSE, then the file will be overwritten.
sep	The character that will separate attributes in a data point.
header	A flag that determines if column names will be output (equivalent to col.names in <a href="#">write.table()</a> ).



row.names	A flag that determines if row names will be output.
close	close stream after writing.
...	Additional parameters that are passed to <code>write.table()</code> .

**Value**

There is no value returned from this operation.

**Author(s)**

Michael Hahsler

**See Also**

[write.table](#)

**Examples**

```
# create data and write 10 points to disk
stream <- DSD_Gaussians(k = 3, d = 5)
stream

write_stream(stream, file="data.txt", n = 10, header = TRUE, info = TRUE)

readLines("data.txt")

# clean up
file.remove("data.txt")

# create a finite stream and write all data to disk using n = -1
stream2 <- DSD_Memory(stream, n = 5)
stream2

write_stream(stream2, file="data.txt", n = -1, header = TRUE, info = TRUE)

readLines("data.txt")

# clean up
file.remove("data.txt")
```

# Index

- \* **DSAggregate**
  - DSAggregate, 8
  - DSAggregate\_Sample, 9
  - DSAggregate\_Window, 10
- \* **DSC\_Macro**
  - DSC\_DBSCAN, 19
  - DSC\_EA, 29
  - DSC\_Hierarchical, 33
  - DSC\_Kmeans, 34
  - DSC\_Macro, 36
  - DSC\_Reachability, 39
  - DSC\_SlidingWindow, 42
- \* **DSC\_Micro**
  - DSC\_BICO, 16
  - DSC\_BIRCH, 18
  - DSC\_DBSTREAM, 21
  - DSC\_DStream, 25
  - DSC\_evoStream, 30
  - DSC\_Micro, 37
  - DSC\_Sample, 41
  - DSC\_Window, 47
- \* **DSC\_TwoStage**
  - DSC\_DBSTREAM, 21
  - DSC\_DStream, 25
  - DSC\_evoStream, 30
  - DSC\_TwoStage, 45
- \* **DSCClassifier**
  - DSCClassifier, 14
  - DSCClassifier\_SlidingWindow, 15
- \* **DSC**
  - animate\_cluster, 4
  - DSC, 12
  - DSC\_Macro, 36
  - DSC\_Micro, 37
  - DSC\_R, 38
  - DSC\_SlidingWindow, 42
  - DSC\_Static, 44
  - DSC\_TwoStage, 45
  - evaluate.DSC, 94
  - get\_assignment, 101
  - plot.DSC, 108
  - predict, 112
  - prune\_clusters, 113
  - read\_saveDSC, 114
  - recluster, 115
- \* **DSD**
  - animate\_data, 6
  - close\_stream, 7
  - DSD, 48
  - DSD\_BarsAndGaussians, 49
  - DSD\_Benchmark, 50
  - DSD\_Cubes, 51
  - DSD\_Gaussians, 52
  - DSD\_Memory, 55
  - DSD\_MG, 57
  - DSD\_Mixture, 59
  - DSD\_mlbenchData, 60
  - DSD\_mlbenchGenerator, 61
  - DSD\_NULL, 62
  - DSD\_ReadDB, 63
  - DSD\_ReadStream, 65
  - DSD\_Target, 68
  - DSD\_UniformNoise, 69
  - DSF, 70
  - get\_points, 103
  - plot.DSD, 110
  - reset\_stream, 116
- \* **DSF**
  - DSF, 70
  - DSF\_Convolve, 73
  - DSF\_Downsampling, 76
  - DSF\_dplyr, 77
  - DSF\_ExponentialMA, 79
  - DSF\_FeatureSelection, 80
  - DSF\_Func, 81
  - DSF\_Scale, 82
- \* **DSOutlier**
  - DSC\_DBSTREAM, 21

- DSC\_DStream, 25
- DSOutlier, 84
- \* **DSRegressor**
  - DSRegressor, 85
  - DSRegressor\_SlidingWindow, 86
- \* **DST**
  - DSAggregate, 8
  - DSC, 12
  - DSCClassifier, 14
  - DSOutlier, 84
  - DSRegressor, 85
  - DST, 87
  - DST\_SlidingWindow, 89
  - DST\_WriteStream, 92
  - evaluate, 93
  - predict, 112
  - stream\_pipeline, 118
  - update, 119
- \* **evaluation**
  - animate\_cluster, 4
  - evaluate, 93
  - evaluate.DSC, 94
- \* **plot**
  - animate\_cluster, 4
  - animate\_data, 6
  - plot.DSC, 108
  - plot.DSD, 110
- add\_cluster (DSD\_MG), 57
- add\_keyframe (MGC), 104
- agreement, 3
- animate (animate\_data), 6
- animate\_cluster, 4, 6, 13, 36, 37, 39, 43, 44, 46, 94, 99, 102, 109, 111, 113–116
- animate\_data, 5, 6, 8, 49–52, 54, 56, 58, 60–62, 64, 67, 69, 70, 72, 104, 109, 111, 117
- animation (animate\_data), 6
- animation::ani.replay(), 5, 6
- BICO (DSC\_BICO), 16
- bico (DSC\_BICO), 16
- BIRCH (DSC\_BIRCH), 18
- birch (DSC\_BIRCH), 18
- change\_alpha (DSC\_DBSTREAM), 21
- close\_stream, 6, 7, 49–52, 54, 56, 58, 60–62, 64, 67, 69, 70, 72, 104, 111, 117
- close\_stream(), 48, 67, 71
- close\_stream.DSD\_ReadCSV (DSD\_ReadStream), 65
- close\_stream.DSD\_ReadDB (DSD\_ReadDB), 63
- close\_stream.DSD\_ReadStream (DSD\_ReadStream), 65
- close\_stream.DSF (DSF), 70
- close\_stream.DST\_WriteStream (DST\_WriteStream), 92
- clue::cl\_agreement(), 4
- cluster::silhouette(), 96, 99
- connection, 66, 67
- D-Stream (DSC\_DStream), 25
- d-stream (DSC\_DStream), 25
- data.frame, 104
- DBI::dbClearResult(), 64
- DBI::dbDisconnect(), 64
- DBI::dbGetQuery(), 64
- DBI::DBI, 64
- DBSCAN (DSC\_DBSCAN), 19
- dbscan (DSC\_DBSCAN), 19
- DBSTREAM, 31
- DBSTREAM (DSC\_DBSTREAM), 21
- dbstream (DSC\_DBSTREAM), 21
- deprecated, 101
- description (DST), 87
- dplyr::dplyr, 78
- dplyr::filter(), 78
- dplyr::mutate(), 78
- dplyr::rename(), 78
- dplyr::select(), 78
- dplyr::transmute(), 78
- DSAggregate, 8, 10, 11, 13, 15, 84, 85, 88, 90, 93, 94, 113, 118, 119
- DSAggregate\_Sample, 9, 9, 11, 41
- DSAggregate\_Window, 9, 10, 10, 90
- DSC, 5, 9, 12, 15, 20, 23, 27, 34–37, 39, 41, 43, 44, 46, 47, 84, 85, 88, 90, 93–95, 99, 101, 102, 108–110, 112–116, 118, 119
- DSC\_BICO, 16, 19, 23, 28, 32, 37, 42, 47
- DSC\_BIRCH, 17, 18, 23, 28, 32, 37, 42, 47
- DSC\_DBSCAN, 19, 30, 34–36, 40, 43
- DSC\_DBSTREAM, 17, 19, 21, 28, 32, 37, 42, 46, 47, 84
- DSC\_DStream, 17, 19, 23, 25, 32, 37, 42, 46, 47, 84
- DSC\_EA, 20, 29, 34–36, 40, 43

- DSC\_evoStream, *17, 19, 23, 28, 29, 30, 37, 42, 46, 47*
- DSC\_Hierarchical, *20, 30, 33, 35, 36, 40, 43*
- DSC\_Kmeans, *20, 30, 34, 34, 36, 40, 42, 43*
- DSC\_Macro, *5, 13, 20, 30, 34, 35, 36, 37, 39, 40, 42–46, 98, 99, 102, 109, 113–116*
- DSC\_Micro, *5, 13, 17, 19, 23, 27, 28, 32, 36, 37, 39, 41–47, 98, 99, 102, 109, 113–116*
- DSC\_R, *5, 13, 20, 23, 27, 34–38, 38, 41, 43, 44, 46, 47, 99, 102, 109, 113–116*
- DSC\_Reachability, *20, 30, 34–36, 39, 43*
- DSC\_Sample, *17, 19, 23, 28, 32, 37, 41, 47*
- DSC\_SlidingWindow, *5, 13, 20, 30, 34–37, 39, 40, 42, 44, 46, 99, 102, 109, 113–116*
- DSC\_Static, *5, 13, 36, 37, 39, 43, 44, 46, 99, 102, 109, 113–116*
- DSC\_TwoStage, *5, 13, 23, 28, 32, 36, 37, 39, 43, 44, 45, 99, 102, 109, 113–116*
- DSC\_Window, *17, 19, 23, 28, 32, 37, 42, 47*
- DSClassifier, *9, 13, 14, 16, 84, 85, 88, 90, 93, 94, 113, 118, 119*
- DSClassifier\_SlidingWindow, *15, 15*
- DSD, *5–8, 27, 46, 48, 50–52, 54, 56, 58–62, 64, 67–72, 74, 76–81, 83, 90, 93–95, 103, 104, 106, 108, 110, 111, 116–119*
- DSD\_BarsAndGaussians, *6, 8, 49, 49, 51, 52, 54, 56, 58, 60–62, 64, 67, 69, 70, 72, 104, 111, 117*
- DSD\_Benchmark, *6, 8, 49, 50, 50, 52, 54, 56, 58, 60–62, 64, 67, 69, 70, 72, 104, 111, 117*
- DSD\_Cubes, *6, 8, 49–51, 51, 54, 56, 58, 60–62, 64, 67, 69, 70, 72, 104, 111, 117*
- DSD\_Gaussians, *6, 8, 49–52, 52, 56, 58, 60–62, 64, 67, 69, 70, 72, 104, 111, 117*
- DSD\_Memory, *6, 8, 49–52, 54, 55, 58, 60–62, 64, 67, 69, 70, 72, 104, 111, 117*
- DSD\_MG, *6, 8, 49–52, 54, 56, 57, 60–62, 64, 67, 69, 70, 72, 104–106, 111, 117*
- DSD\_Mixture, *6, 8, 49–52, 54, 56, 58, 59, 61, 62, 64, 67, 69, 70, 72, 104, 111, 117*
- DSD\_mlbenchData, *6, 8, 49–52, 54, 56, 58, 60, 60, 62, 64, 67, 69, 70, 72, 104, 111, 117*
- DSD\_mlbenchGenerator, *6, 8, 49–52, 54, 56, 58, 60, 61, 61, 62, 64, 67, 69, 70, 72, 104, 111, 117*
- DSD\_NULL, *6, 8, 49–52, 54, 56, 58, 60–62, 62, 64, 67, 69, 70, 72, 104, 111, 117*
- DSD\_R, *52, 54, 56, 59, 61, 64, 67*
- DSD\_R (DSD), *48*
- DSD\_ReadCSV, *7, 117*
- DSD\_ReadCSV (DSD\_ReadStream), *65*
- DSD\_ReadDB, *6–8, 49–52, 54, 56, 58, 60–62, 63, 67, 69, 70, 72, 104, 111, 117*
- DSD\_ReadStream, *6–8, 49–52, 54, 56, 58, 60–62, 64, 65, 69, 70, 72, 104, 111, 117*
- DSD\_ScaleStream, *68*
- DSD\_Target, *6, 8, 49–52, 54, 56, 58, 60–62, 64, 67, 68, 70, 72, 104, 111, 117*
- DSD\_UniformNoise, *6, 8, 49–52, 54, 56, 58, 60–62, 64, 67, 69, 69, 72, 104, 111, 117*
- DSF, *6–8, 49–52, 54, 56, 58, 60–62, 64, 67–70, 70, 74–76, 78–81, 83, 104, 111, 117, 118*
- DSF\_Convolve, *72, 73, 76, 78–81, 83, 117*
- DSF\_Downsampling, *72, 75, 76, 78–81, 83*
- DSF\_dplyr, *72, 75, 76, 77, 79–81, 83*
- DSF\_ExponentialMA, *72, 75, 76, 78, 79, 80, 81, 83*
- DSF\_FeatureSelection, *72, 75, 76, 78, 79, 80, 81, 83*
- DSF\_Func, *72, 75, 76, 78–80, 81, 83*
- DSF\_Scale, *68, 72, 75, 76, 78–81, 82*
- DSF\_Scale(), *68*
- DSFP, *72, 88*
- DSOutlier, *9, 13, 15, 23, 28, 84, 85, 88, 90, 93, 94, 113, 118, 119*
- DSOutlier\_DBSTREAM (DSC\_DBSTREAM), *21*
- DSOutlier\_DStream (DSC\_DStream), *25*
- DSRegressor, *9, 13, 15, 84, 85, 86, 88, 90, 93, 94, 113, 118, 119*
- DSRegressor\_SlidingWindow, *85, 86*
- DST, *8, 9, 13, 15, 37, 71, 73, 84, 85, 87, 88, 90, 93, 94, 103, 112, 113, 118, 119*
- DST\_Multi, *88, 118*
- DST\_Runner (stream\_pipeline), *118*
- DST\_SlidingWindow, *9, 13, 15, 43, 84–86, 88, 89, 93, 94, 113, 118, 119*
- DST\_WriteStream, *7, 9, 13, 15, 84, 85, 88, 90, 92, 94, 113, 118, 119*

- dstream (DSC\_DStream), 25
- evaluate, 5, 9, 13, 15, 84, 85, 88, 90, 93, 93, 99, 113, 118, 119
- evaluate.DSC, 5, 13, 36, 37, 39, 43, 44, 46, 93, 94, 94, 102, 109, 113–116
- evaluate\_static (evaluate), 93
- evaluate\_static(), 94
- evaluate\_static.DSC (evaluate.DSC), 94
- evaluate\_stream (evaluate), 93
- evaluate\_stream(), 5, 94
- evaluate\_stream.DSC (evaluate.DSC), 94
- filter\_diff (DSF\_Convolve), 73
- filter\_Hamming (DSF\_Convolve), 73
- filter\_MA (DSF\_Convolve), 73
- filter\_Sinc (DSF\_Convolve), 73
- formula, 17–19, 22, 26, 29, 31, 33, 35, 39
- fpc::cluster.stats(), 96, 99
- get\_assignment, 5, 13, 36, 37, 39, 43, 44, 46, 99, 101, 109, 113–116
- get\_attraction (DSC\_DStream), 25
- get\_centers (DSC), 12
- get\_centers(), 27, 46
- get\_clusters (DSD\_MG), 57
- get\_clusters(), 37
- get\_copy (DSC), 12
- get\_keyframes (MGC), 104
- get\_macroclusters (DSC), 12
- get\_macroweights (DSC), 12
- get\_microclusters (DSC), 12
- get\_microweights (DSC), 12
- get\_model (DST), 87
- get\_points, 6, 8, 49–52, 54, 56, 58, 60–62, 64, 67, 69, 70, 72, 103, 111, 117
- get\_points(), 8, 48, 62, 66, 71, 81
- get\_points.DSAggregate (DSAggregate), 8
- get\_points.DSF (DSF), 70
- get\_shared\_density (DSC\_DBSTREAM), 21
- get\_weights (DSC), 12
- get\_weights(), 8, 46
- get\_weights.DSAggregate (DSAggregate), 8
- graphics::pairs(), 109, 110
- graphics::plot.default(), 109, 110
- hclust(), 33
- keyframe (MGC), 104
- lm(), 86
- magrittr::%>%, 71
- MGC, 57, 58, 104, 105
- MGC\_Function (MGC), 104
- MGC\_Linear (MGC), 104
- MGC\_Noise (MGC), 104
- MGC\_Random (MGC), 104
- MGC\_Static (MGC), 104
- microToMacro (DSC\_Macro), 36
- microToMacro(), 36
- nclusters (DSC), 12
- plot (plot.DSD), 110
- plot(), 27, 48, 84
- plot.DSC, 5, 6, 13, 36, 37, 39, 43, 44, 46, 99, 102, 108, 111, 113–116
- plot.DSC(), 13, 110
- plot.DSC\_DBSTREAM (DSC\_DBSTREAM), 21
- plot.DSC\_DStream (DSC\_DStream), 25
- plot.DSD, 5, 6, 8, 49–52, 54, 56, 58, 60–62, 64, 67, 69, 70, 72, 104, 109, 110, 117
- plot.DSD(), 108
- pow2 (DSF\_Convolve), 73
- predict, 5, 9, 13, 15, 36, 37, 39, 43, 44, 46, 84, 85, 88, 90, 93, 94, 99, 102, 109, 112, 114–116, 118, 119
- predict(), 13, 23, 88, 95, 101
- predict.DST\_SlidingWindow (DST\_SlidingWindow), 89
- prune\_clusters, 5, 13, 36, 37, 39, 43, 44, 46, 99, 102, 109, 113, 113, 115, 116
- read.table(), 66, 67
- read\_saveDSC, 5, 13, 36, 37, 39, 43, 44, 46, 99, 102, 109, 113, 114, 114, 116
- readDSC (read\_saveDSC), 114
- readDSC(), 13
- readLines(), 66, 67
- readRDS, 115
- recluster, 5, 13, 36, 37, 39, 43, 44, 46, 99, 102, 109, 113–115, 115
- recluster(), 20, 32, 33, 35, 36, 40
- remove\_cluster (DSD\_MG), 57
- remove\_info (get\_points), 103
- remove\_keyframe (MGC), 104
- reset\_stream, 6, 8, 49–52, 54, 56, 58, 60–62, 64, 67, 69, 70, 72, 104, 111, 116

`reset_stream()`, [48](#), [60](#), [64](#), [67](#), [71](#), [103](#)  
`reset_stream.DSF (DSF)`, [70](#)  
`rpart()`, [15](#)

`saveDSC (read_saveDSC)`, [114](#)  
`saveDSC()`, [13](#)  
`saveRDS`, [115](#)  
`scale`, [83](#)  
`seek()`, [117](#)  
`Shape_Block (MGC)`, [104](#)  
`Shape_Gaussian (MGC)`, [104](#)  
`stats::filter`, [75](#)  
`stats::kmeans()`, [35](#)  
`stream_pipeline`, [9](#), [13](#), [15](#), [84](#), [85](#), [88](#), [90](#),  
[93](#), [94](#), [113](#), [118](#), [119](#)

`update`, [9](#), [13](#), [15](#), [84](#), [85](#), [88](#), [90](#), [93](#), [94](#), [113](#),  
[118](#), [119](#)  
`update()`, [13](#), [20](#), [23](#), [33](#), [35](#), [37](#), [40](#), [71](#), [88](#),  
[118](#)  
`update.DSAggregate (DSAggregate)`, [8](#)  
`update.DSC_R (DSC_R)`, [38](#)  
`update.DSF (DSF)`, [70](#)  
`update.DST_SlidingWindow  
(DST_SlidingWindow)`, [89](#)

`write.table`, [121](#)  
`write.table()`, [120](#), [121](#)  
`write_stream`, [120](#)  
`write_stream()`, [92](#)