

Package ‘ASW’

October 31, 2023

Title Clustering Algorithms for Optimizing the Average Silhouette Width

Version 0.0.1

Author Minh Long Nguyen <edelweiss611428@gmail.com>

Maintainer Minh Long Nguyen <edelweiss611428@gmail.com>

Description This package implements clustering algorithms for optimizing the Average Silhouette Width, including PAMSil, Efficient Optimum Silhouette (effOSil), and Scalable Optimum Silhouette (scalOSil).

License GPL (>= 3)

Encoding UTF-8

Roxygen list(markdown = TRUE)

RoxygenNote 7.2.3

LinkingTo Rcpp

Imports Rcpp,
cluster,
stats

R topics documented:

effOSil	1
Init	3
PAMSil	4
scalOSil	5
Silhouette	6
Index	7

effOSil	<i>The Efficient Optimum Silhouette algorithm</i>
---------	---------------------------------------------------

Description

This function implements the Efficient Optimum Silhouette (effOSil) algorithm.

Usage

```
effOSil(dx, K, initMethod, variant)
```

Arguments

<code>dx</code>	A "dist" object, which can be computed using <code>stats::dist()</code> .
<code>K</code>	An integer vector (or scalar) specifying the numbers of clusters. By default, <code>K = 2:12</code> .
<code>initMethod</code>	A character vector (or string) specifying initialization methods. By default, <code>initMethod = "average"</code> . See <code>?Init</code> for more details.
<code>variant</code>	A character string specifying a variant. Options include "efficient" and "original". If <code>variant = "original"</code> , the original OSil algorithm is used. If <code>variant = "efficient"</code> , effOSil is used. By default, <code>variant = "efficient"</code> .

Details

This function implements the Efficient Optimum Silhouette (effOSil) algorithm, an $O(N)$ runtime improvement of the original, computationally expensive OSil algorithm proposed by Batool & Hennig (2021) (N is the number of observations). An implementation of the original OSil algorithm is also available for run time comparisons.

Value

best_clustering The effOSil clustering achieving the highest ASW value.

best_asw The highest ASW value.

k The estimated number of clusters.

clusterings The effOSil clustering solutions for all k in K .

asw The ASW values associated with the effOSil clusterings.

nIter The numbers of iterations needed for convergence.

Author(s)

Minh Long Nguyen <edelweiss611428@gmail.com>

References

Batool, F. and Hennig, C., 2021. Clustering with the average silhouette width. Computational Statistics & Data Analysis, 158, p.107190.

Examples

```
dx = dist(faithful)
effC = effOSil(dx, 2:8)
par(mfrow = c(2,1))
plot(faithful, col = effC$best_clustering, pch = 4)
plot(2:8, effC$asw, xlab = "k", ylab = "ASW")
```

Init

*Initialization methods for the Optimum Silhouette algorithm***Description**

This function computes an initialization for the Optimum Silhouette algorithm.

Usage

```
Init(dx, k, initMethod)
```

Arguments

<code>dx</code>	dx A "dist" object, which can be computed using <code>stats::dist()</code> .
<code>k</code>	An integer scalar specifying the number of clusters.
<code>initMethod</code>	A character vector (or string) specifying initialization methods. Options include any combination of "pam", "average", "single", "complete", "ward.D", "ward.D2", "mcquitty", "median", and "centroid". By default, <code>initMethod = "average"</code> .

Details

This function computes an initialization for the Optimum Silhouette algorithm, but it can be used as a stand-alone clustering method.

Value

clustering An initialized clustering.
asw The ASW associated with the initialized clustering.
method The "best" initialization method.

Author(s)

Minh Long Nguyen <edelweiss611428@gmail.com>

References

Batool, F. and Hennig, C., 2021. Clustering with the average silhouette width. *Computational Statistics & Data Analysis*, 158, p.107190. Batool, F., 2019. Initialization methods for optimum average silhouette width clustering. *arXiv preprint arXiv:1910.08644*.

Examples

```
x = faithful
dx = dist(x)
Initres = Init(dx, 2, c("pam", "average", "complete"))
plot(x, col = Initres$clustering, pch = 4)
print(paste(Initres$method, "achieves the highest ASW value"))
```

Description

This function implements the PAMSil algorithm.

Usage

```
PAMSil(dx, K)
```

Arguments

dx	A "dist" object, which can be computed using <code>stats::dist()</code> .
K	An integer vector (or scalar) specifying the numbers of clusters. By default, <code>K = 2:12</code> .

Details

This function implements the PAMSil algorithm proposed by Van der Laan et al. (2003). It is a k-medoids clustering algorithm whose objective function is the Average Silhouette Width.

Value

best_clustering The PAMSil clustering achieving the highest ASW value.
best_asw The highest ASW value.
best_medoids The medoids associated with the clustering maximize the ASW.
k The estimated number of clusters.
clusterings The PAMSil clustering solutions for all `k` in `K`.
asw The ASW values associated with the PAMSil clusterings.
medoids The medoids associated with the clustering solutions.
nIter The numbers of iterations needed for convergence.

Author(s)

Minh Long Nguyen <edelweiss611428@gmail.com>

References

Van der Laan, M., Pollard, K. and Bryan, J., 2003. A new partitioning around medoids algorithm. *Journal of Statistical Computation and Simulation*, 73(8), pp.575-584.

Examples

```
dx = dist(faithful)
pamsilC = PAMSil(dx, 2:8)
par(mfrow = c(2,1))
plot(faithful, col = pamsilC$best_clustering, pch = 4)
plot(2:8, pamsilC$asw, xlab = "k", ylab = "ASW")
```

scalOSil

*The Scalable Optimum Silhouette algorithm***Description**

This function implements the Scalable Optimum Silhouette algorithm.

Usage

```
scalOSil(dx, K, n, ns, rep, initMethod, variant)
```

Arguments

<code>dx</code>	A "dist" object, which can be computed using <code>stats::dist()</code> .
<code>K</code>	An integer vector (or scalar) specifying the numbers of clusters. By default, <code>K = 2:12</code> .
<code>n</code>	An integer specifying the sample size. If not specified (NULL), <code>n</code> is set to $0.2 \cdot N$ where <code>N</code> is the number of observations.
<code>ns</code>	An integer specifying the number of random samples used in each instance. By default, <code>ns = 1</code> .
<code>rep</code>	An integer specifying the number of scalOSil instances. By default, <code>rep = 10</code> .
<code>initMethod</code>	A character vector (or string) specifying initialization methods. By default, <code>initMethod = "average"</code> . See <code>?Init</code> for more details.
<code>variant</code>	A character string specifying a variant. Options include "scalable" and "original". If <code>variant = "original"</code> , the original FOSil algorithm is used. If <code>variant = "scalable"</code> , scalOSil is used. By default, <code>variant = "scalable"</code> .

Details

This function implements the Scalable Optimum Silhouette (scalOSil) algorithm, an $O(n)$ runtime improvement of the original, computationally expensive Fast OSil (FOSil) algorithm proposed by Batool & Hennig (2021) (`n` is the sample size). An implementation of the original FOSil algorithm is also available for run time comparisons.

Value

best_clustering The scalOSil clustering achieving the highest ASW value.
best_asw The highest ASW value.
k The estimated number of clusters.
clusterings The scalOSil clustering solutions for all `k` in `K`.
asw The ASW values associated with the scalOSil clusterings.

Author(s)

Minh Long Nguyen <edelweiss611428@gmail.com>

References

Batool, F. and Hennig, C., 2021. Clustering with the average silhouette width. Computational Statistics & Data Analysis, 158, p.107190.

Examples

```
dx = dist(faithful)
scalC = scal0Sil(dx, 2:8)
par(mfrow = c(2,1))
plot(faithful, col = scalC$best_clustering, pch = 4)
plot(2:8, scalC$asw, xlab = "k", ylab = "ASW")
```

Silhouette	<i>Silhouette Width</i>
------------	-------------------------

Description

This function computes the Silhouette Widths for all data points in the dataset.

Usage

```
Silhouette(C, dx)
```

Arguments

C An integer vector specifying a k-partition of the dataset. $\min(C)$ must be 1 and $\max(C)$ must be k.

dx A "dist" object, which can be computed using `stats::dist()`.

Value

A numeric matrix of class "silhouette" containing three columns

cluster A clustering of the dataset.

neighbor The clustering labels of the nearest clusters for all data points.

sil_width The silhouette widths of data points.

Author(s)

Minh Long Nguyen <edelweiss611428@gmail.com>

References

Rousseeuw, P.J. (1987) Silhouettes: A graphical aid to the interpretation and validation of cluster analysis. J. Comput. Appl. Math., 20, 53–65.

Examples

```
library("cluster")
dx = dist(faithful)
C = pam(dx, 2)$clustering
plot(Silhouette(C,dx))
```

Index

eff0Sil, [1](#)

Init, [3](#)

PAMSil, [4](#)

scal0Sil, [5](#)

Silhouette, [6](#)