

# Package ‘HCsnip’

October 9, 2013

**Type** Package

**Title** Semi-supervised adaptive-height snipping of the Hierarchical Clustering tree

**Version** 1.0.0

**Author** Askar Obulkasim

**Maintainer** Askar Obulkasim <askar.wubulikasimu@vumc.nl>

**Description** Decompose given hierarchical clustering tree into non-overlapping clusters in a semi-supervised way by using available patients follow-up information as guidance. Contains functions for snipping HC tree, various cluster quality evaluation criteria, assigning new patients to one of the two given HC trees, testing the significance of clusters with permutation argument and clusters visualization using sample’s molecular entropy.

**License** GPL (>= 2)

**LazyLoad** yes

## **Depends**

R(>= 2.10.0), survival, coin, fpc, clusterRepro, impute,randomForestSRC, sm, sigaR, Biobase

**biocViews** Microarray, Bioinformatics, aCGH, GeneExpression, Clustering

## **R topics documented:**

|                             |    |
|-----------------------------|----|
| HCsnip-package . . . . .    | 2  |
| BullingerLeukemia . . . . . | 2  |
| cluster_pred . . . . .      | 3  |
| EnvioPlot . . . . .         | 5  |
| HCsnipper . . . . .         | 6  |
| measure . . . . .           | 8  |
| perm_test . . . . .         | 9  |
| RSF_eval . . . . .          | 12 |
| surv_measure . . . . .      | 13 |
| TcgaGBM . . . . .           | 14 |
| TwoHC_assign . . . . .      | 15 |
| TwoHC_perm . . . . .        | 17 |

**Index**

20

---

|                |   |
|----------------|---|
| HCsnip-package | <i>Semi-supervised adaptive-height snipping of the Hierarchical Clustering tree</i> |
|----------------|---|

---

**Description**

This package contains functions for extracting meaningful clusters from a HC tree. Rather than cutting the HC tree at a fixed height (as existing methods do), it snips the tree at variable heights to extract hidden clusters. Cluster extraction process uses both the data type from which HC tree is derived and the available patients follow-up information. Functions for testing the significance of extracted clusters and cluster visualization using sample's molecular entropy are also given. If two HC trees are presented which may correspond to the two treatment groups, this package also includes functions for optimally assigning new patients to one of the two HC trees and calculate the expected gain in terms of follow-up.

**Details**

Package: HCsnip  
Type: Package  
License: GPL (>= 2)

**Author(s)**

Askar Obulkasim Maintainer: Askar Obulkasim <askar.wubulikasimu@vumc.nl>

**References**

Obulkasim, A. et al., (2013). "Semi-supervised adaptive-height snipping of the Hierarchical Clustering tree", submitted.

---

|                   |                      |
|-------------------|----------------------|
| BullingerLeukemia | <i>Leukemia data</i> |
|-------------------|----------------------|

---

**Description**

Gene expression profiles of adult acute myeloid leukemia patients. Contains 116 samples and 1571 genes. The follow-up data also included.

**Usage**

data(BullingerLeukemia)

**Format**

An object of list class.

**Note**

The complete dataset is available at the Gene Expression Omnibus [www.ncbi.nlm.nih.gov/geo/](http://www.ncbi.nlm.nih.gov/geo/), accession number GSE425.

**Source**

Bullinger,L. et al., (2004). "Use of gene-expression profiling to identify prognostic subclasses in adult acute myeloid leukemia", *N Engl J Med.*, 350, 1605-1616.

**Examples**

```
data(BullingerLeukemia)
names(BullingerLeukemia)
```

---

|              |                                   |
|--------------|-----------------------------------|
| cluster_pred | <i>Semi-supervised clustering</i> |
|--------------|-----------------------------------|

---

**Description**

For a given partition, this function assigns new samples to one of the clusters in the partition. Partition (clustering) is composed of non-overlapping clusters

**Usage**

```
cluster_pred(X, partition, surv.time = NULL, status = NULL, te.index, minclus = 4,
             te.surv.time = NULL, te.status = NULL, method = "conc", maxmiss = 30,
             plot.it = FALSE, ...)
```

**Arguments**

|           |   |
|-----------|---|
| X         | An object of class <a href="#">ExpressionSet</a> or data matrix in which columns are assumed to represent the samples, and rows represents the sample's features. X can also be a square distance matrix or object class of <a href="#">dist</a> . It must include the data set from which <i>partition</i> is obtained and the data set of test samples. |
| partition | A numeric vector contains the non-overlapping cluster labels.   |
| surv.time | A numeric vector contains follow-up time of patients in the partition.  |
| status    | A binary vector contains survival status of patients in the partition, 0 = alive, 1 = dead.   |
| te.index  | A numeric vector contains the indices of columns in X corresponds to the test samples.  |
| minclus   | The minimum number samples allowed to form a cluster for the test set. This is to avoid returning tiny clusters and reduce the effect of outliers.  |

|                           |   |
|---------------------------|---|
| <code>te.surv.time</code> | An optional vector contains follow-up time of patients in the test set. If supplied with <code>te.status</code> , the logrank test <i>p-value</i> is calculated for the test set.     |
| <code>te.status</code>    | An optional vector contains survival status of patients in the test set. If supplied with <code>te.surv.time</code> , the logrank test <i>p-value</i> is calculated for the test set. |
| <code>method</code>       | Type of methods to use in assigning test samples to one of the clusters in the partition. Must be either Ward distance 'ward' or Harrel's concordance index 'conc' (default).         |
| <code>maxmiss</code>      | Maximum percentage of missing values per row in <i>X</i>  |
| <code>plot.it</code>      | If TRUE and follow-up data of the test samples are given, Kaplan Meier curve(s) will be generated for each cluster in the test set.   |
| <code>...</code>          | Arguments for <code>impute.knn</code> from the <b>impute</b> package.   |

### Details

User has two options to assign test set to one of the clusters in the partition. One option is to use the Ward distance. Specifically, an average distance is calculated between a test sample and samples in each cluster in the partition, separately. The test sample is assigned to a cluster for which average distance is the smallest. Follow-up data is not required for this option.

Second option is to use the Harrel's concordance index (Harrel et al., 1982). For this option both main and follow-up data corresponds to the given partition are required. Main data is used to find the pseudo nearest neighbours (PNN) of a test sample (Obulkasim et a., 2011), and follow-up data is used to check how much PNN's follow-up info is concordant with follow-up info of samples in each cluster. The test sample is assigned to a cluster for which average concordance is the highest.

Before selecting either one of the options, we recommend user to check the correlation between main data and follow-up info (e.g. using global test). If correlation is relatively large, we recommend to use 'conc' option, and vice versa.

### Value

If `plot.it` is FALSE, function returns a vector of cluster labels of the test set. If TRUE and follow-up data of the test set are given, function returns a list object contains following components:

|                    |  |
|--------------------|--|
| <code>St</code>    | a data frame with following five columns:                                |
|                    | [,1] the unique survival times found in the test set's follow-up.        |
|                    | [,2] the cluster labels of the test set.                                 |
|                    | [,3] the number of patients at risk at each unique time point.           |
|                    | [,4] the survival probability at each unique time point.                 |
|                    | [,5] the variance of the survival probability at each unique time point. |
| <code>value</code> | logrank test <i>p-value</i> for the test set                             |

### Author(s)

Askar Obulkasim

## References

- Obulkasim,A. et al., (2013). "Semi-supervised adaptive-height snipping of the Hierarchical Clustering tree", submitted.
- Harrel,E.F. et al., (1982). "Evaluating the yield of medical tests", *JAMA*, 247, 2543-2546.
- Obulkasim,A. et al., (2011). "Stepwise classification of cancer samples using clinical and molecular data", *BMC Bioinformatics*, 12, 422.
- Troyanskaya,O. et al., (2001). "Missing value estimation methods for DNA microarrays". *Bioinformatics*, 17, 520-525.

## See Also

[TwoHC\\_assign](#)

## Examples

```
data(BullingerLeukemia)
em <- BullingerLeukemia[["em"]]
st <- BullingerLeukemia[["surv.time"]]
event <- BullingerLeukemia[["status"]]
cl <- HCsnipper(em[, 1:60], min = 5)
result <- cluster_pred(X = em[, 1:100], partition = cl$partitions[1, ], surv.time = st[1:60],
                      status = event[1:60], te.index = 61:100)
names(result)
```

---

EnvioPlot

*Visualize cluster's molecular entropy by violin plot*

---

## Description

This function first calculates the entropy of each sample in the given partition, and make a violin plot for each cluster. If clusters are different in term of their molecular profiles, then one may expect density differences in the violin plot.

## Usage

```
EnvioPlot(X, method = "knn", parti, horizontal = FALSE,
          col = NULL, names = NULL, ...)
```

## Arguments

- |                     |   |
|---------------------|---|
| <code>X</code>      | An object of class <a href="#">ExpressionSet</a> or data matrix from which partition is obtained. Columns are assumed to represent samples, and rows represent sample's features. Missing values are allowed. |
| <code>method</code> | Type of method to calculate sample's molecular entropy. Either <i>knn</i> or <i>normal</i> .  |
| <code>parti</code>  | A partition for which violin plot to be made.   |

|                         |  |
|-------------------------|--|
| <code>horizontal</code> | Should boxes are organized horizontally? default is FALSE.   |
| <code>col</code>        | A vector of colors for each cluster. Should be equal the number of clusters in $x$ .                   |
| <code>names</code>      | A vector of labels for each cluster.   |
| <code>...</code>        | Arguments for <code>hdEntropy</code> to calculate the molecular entropy from the <b>sigar</b> package. |

### Value

The entropy estimate is returned as a numeric.

### Author(s)

Askar Obulkasim

### References

van Wieringen, N.W. and van der Vaart, W.A. (2010) "Statistical analysis of the cancer cell's molecular entropy using high-throughput data", *Bioinformatics*, 27, 556-563.

Adler, D. (2005). `vioplot`: A violin plot is a combination of a box plot and a kernel density plot, R package, <http://cran.r-project.org/web/packages/vioplot/index.html>.

### See Also

[hdEntropy](#)

### Examples

```
data(BullingerLeukemia)
em <- BullingerLeukemia[["em"]]
cl <- cutree(hclust(as.dist(1 - cor(em[, 1:60])), method = "ward"), k = 2)
result <- EnvioPlot(X = em[, 1:60], parti = cl)
```

---

HCsnipper

*HC tree snipper*

---

### Description

This function snips given hierarchical clustering (HC) at variable heights to extract all possible partitions. Each partition (clustering) is composed of non-overlapping clusters.

### Usage

```
HCsnipper(X, hc = NULL, dis = NULL, dis.method = "cor", link.method = "ward",
          minclus = 4, maxmiss = 30, ...)
```

**Arguments**

|                          |  |
|--------------------------|--|
| <code>X</code>           | An object of class <code>ExpressionSet</code> or data matrix from which HC tree will be derived. Columns are assumed to represent the samples, and rows represent the sample's features (genes). Missing values are allowed. |
| <code>hc</code>          | HC tree from which partitions to be extracted. Must be an object class of <code>hclust</code> . This is an optional argument, but if given <code>X</code> and <code>dis</code> will be ignored.                              |
| <code>dis</code>         | A square distance matrix or object class of <code>dist</code> from which HC tree to be derived. This is an optional argument, if given <code>X</code> will be ignored.   |
| <code>dis.method</code>  | The distance measure to be used. This must be one of the methods acceptable for <code>dist</code> function or the Pearson correlation 'cor' (default).   |
| <code>link.method</code> | The agglomeration method to be used. This should be one of "ward" (default), "single", "complete", "average", "mcquitty", "median" or "centroid".  |
| <code>minclus</code>     | The minimum number of samples allowed to form a cluster. This parameter is inversely proportional to the number of partitions returned. e.g. large values returns less number clusters, and vice versa.                      |
| <code>maxmiss</code>     | Maximum percentage of missing values per row in <code>X</code>   |
| <code>...</code>         | Arguments for <code>impute.knn</code> from the <code>impute</code> package for missing values imputation in <code>X</code> .   |

**Details**

For given HC tree, this function snips it at all possible places to extract partitions under the following conditions:

- Singleton is not allowed.
- Snipping places are chosen so that only the samples which are neighbours in the leaf node ordering (see `order(hc)`) are allowed to form a cluster.

The last constraint guarantees that sniping does not change the HC tree structure considerably. For example, samples located in far left in the HC tree will not be joined with samples located in far right. The number of partitions return by function depends not only on the `minclus` argument, but also the shape of the HC tree. Large number of partitions can be returned from a balanced HC tree than a skewed one.

**Value**

This function returns a list object contains following objects:

|                         |  |
|-------------------------|--|
| <code>partitions</code> | a matrix in which rows represent partitions and columns represent samples.                               |
| <code>hc</code>         | HC tree from which partitions are extracted.   |
| <code>dat</code>        | data matrix. If <code>X</code> has missing values, this will be missing values imputed full data matrix. |
| <code>dis</code>        | the distance matrix used   |
| <code>dis.m</code>      | the distance measure used  |
| <code>link.m</code>     | the agglomeration method used  |

**Author(s)**

Askar Obulkasim

**References**

Obulkasim,A. et al., (2013). "Semi-supervised adaptive-height snipping of the Hierarchical Clustering tree", submitted.

Troyanskaya,O. et al., (2001). "Missing value estimation methods for DNA microarrays". *Bioinformatics*, 17, 520-525.

**Examples**

```
data(BullingerLeukemia)
em <- BullingerLeukemia[["em"]]
c1 <- HCsnipper(em[, 1:60], minclus = 5)
## To check returned partitions size
table(apply(c1[["partitions"]], 1, function(x) length(unique(x))))
```

---

measure

*Evaluate cluster quality*

---

**Description**

Function to evaluate the overall quality of a partition (composed of non-overlapping clusters) by user defined criteria.

**Usage**

```
measure(parti, dis, X = NULL, method = "g2", maxmiss = 30, ...)
```

**Arguments**

|                      |   |
|----------------------|---|
| <code>parti</code>   | Partition to be evaluated.  |
| <code>dis</code>     | A square distance matrix or class object of <code>dist</code> corresponding to $x$ .  |
| <code>X</code>       | data matrix corresponding to the <i>parti</i> . Columns are assumed to represent the samples, and rows represent the sample's features. Missing values are allowed. This is an optional argument, but If type is set to 'igp', then matrix must be given. |
| <code>method</code>  | Type of evaluation measure to use for assessing the quality of clusters in $x$ . Default is Goodman and Kruskal index $g2$ .  |
| <code>maxmiss</code> | Maximum percentage of missing values per row in <i>dat</i>  |
| <code>...</code>     | Arguments for function <code>cluster.stats</code> from the <b>fpc</b> package. See details below.   |

## Details

Numerous cluster quality measuring criteria have been proposed. This package includes only a few well known ones. Except for the 'c.index' and the in group proportion 'igp', rest of the criteria come from the function `cluster.stats` in **fpc** package. For latter one, please see the returned arguments of the `cluster.stats` function before you decide which criteria to choose. Note that, the value returned by different criteria has different meaning. For example. the larger the Goodman and Kruskal index 'g2' the better, for the index G3 'g3' the smaller the better. Thus, interpret returned value accordingly.

## Value

A numeric value representing the quality of partition under consideration.

## Author(s)

Askar Obulkasim

## References

Hennig,C. (2010). fpc: Flexible procedures for clustering, R package, <http://CRAN.R-project.org/package=fpc>.

Kapp,A.V. and Tibshirani,R. (2007) "Are clusters found in one dataset present in another dataset?", *Biostatistics*, 8, 9-31.

Obulkasim,A. et al., (2013). "Semi-supervised adaptive-height snipping of the Hierarchical Clustering tree", submitted.

## See Also

[surv\\_measure](#)

## Examples

```
data(BullingerLeukemia)
em <- BullingerLeukemia[["em"]]
cl <- HCsnipper(em[, 1:60], minclus = 5)
m <- apply(cl[["partitions"]], 1, function(x) measure(parti = x, dis = 1 - cor(em[, 1:60])))
```

---

perm\_test

*A function to select an optimal partition (clustering) from large number of candidates and calculate the p-value for it.*

---

## Description

For a given set of partitions (each partition is composed of non-overlapping clusters), this function uses two types of data to evaluate each partition and select the optimal one which has the highest rank in terms of both data type (presumed that *score1* and *score2* were from two different data source). Permutation approach used to calculate the corrected *p-value* of the selected partition.

## Usage

```
perm_test(partitions, surv.time, status, score1 = NULL, score2, method = "BIC", nperm = 1000)
```

## Arguments

|            |   |
|------------|---|
| partitions | A matrix in which rows represent partitions and columns represent samples   |
| surv.time  | A numeric vector contains follow-up time of patients in the partition   |
| status     | A binary vector contains survival status of patients in the partition, 0 = alive, 1 = dead  |
| score1     | A numeric vector contains the quality score for each partition. Scores are assumed to be calculated using the follow-up data. Note, prepare this vector in a way that high value corresponds to good quality partition.                       |
| score2     | A numeric vector contains the quality score for each partition calculated by using any data type except for follow-up. The same as <i>score1</i> this vector must be prepared in a way that high value corresponds to good quality partition. |
| method     | Type of partition evaluation measure to use. Must be the same as the type of measure used in calculating the <i>score1</i> . Default is 'BIC'   |
| nperm      | The number of permutations.   |

## Details

When studying association of cluster membership with follow-up data, we cannot use the standard testing procedures. Because *score1* is already used the follow-up data. Thus, we would use the follow-up data twice and the resulting *p-value* is likely to be too small. We avoid this bias by also applying the semi-supervised partition selection under the null-hypothesis. This null-hypothesis is simply the absence of association between the data type used to generate the *score2* and the follow-up. Our partition selection in combination with a suitable test statistic is designed to detect associations that can be represented by groups of samples. We adapt the *p-value* computation as follows:

1. Use a suitable test statistic (e.g. log-rank for time-to-event data and chi-square for nominal data) to compute the conditional *p-value* given the cluster labels in the selected partition:  $p_{obs}$ .
2. For  $i = 1 \dots nperm$ :
  - (a) Randomly permute follow-up data among the samples.
  - (b) Apply exactly the same type of evaluation measure to evaluate all partitions, e.g. generate new *score1*, but *score2* is fixed. Selected the best partition as before.
  - (c) Conditional on the resulting partition, compute *p-value*  $p_i$ .

3. Finally, the  $p$ -value of interest is equal the number of time  $p_i$  smaller (or equal) than the  $p_{obs}$  divided by the number of permutations ran.

Here,  $p$  satisfies a crucial property of  $p$ -value: it is uniformly distributed when the null-hypothesis is true, because then  $p_{obs}$  and  $p_i$  are exchangeable random variables. The exchangeability is a result from the null-hypothesis and the use of exactly the same procedures to compute  $p_{obs}$  and  $p_i$ .

### Value

A list object contains following objects:

|        |  |
|--------|--|
| obs.p  | Observed $p$ -value                        |
| perm.p | A vector of $p$ -values from permutations. |
| best   | Selected optimal partition                 |

### Author(s)

Askar Obulkasim

### References

Obulkasim,A. et al., (2013). "Semi-supervised adaptive-height snipping of the Hierarchical Clustering tree", submitted.

### See Also

[TwoHC\\_perm](#)

### Examples

```
data(BullingerLeukemia)
em <- BullingerLeukemia[["em"]]
st <- BullingerLeukemia[["surv.time"]]
event <- BullingerLeukemia[["status"]]
cl <- HCsnipper(em[, 1:60], min = 5)
m <- apply(cl[["partitions"]], 1, function(x) measure(parti = x, dis = 1-cor(em[, 1:60])))
s <- apply(cl[["partitions"]], 1, function(x) surv_measure(x, st[1:60], event[1:60]))
result <- perm_test(cl[["partitions"]], st[1:60], event[1:60], score1 = s, score2 = m, nperm = 10)

### Visualize cluster differences in terms of Entropy.
H <- EnvioPlot(X = em[, 1:60], parti = result[["best"]])
```

RSF\_eval

*Function to calculate error rate using the Random Survival Forest***Description**

This function constructs survival forest using training sample's follow-up as response and cluster labels as covariate. Constructed forest is used to calculate cumulative hazard function (CHF) for each new sample based its cluster label. CHFs are compared with new samples' actual survival time to calculate the error rate. Error rate ranges from 0 to 1, with 0 representing perfect.

**Usage**

```
RSF_eval(partition, surv.time, status, te.partition, te.surv.time, te.status, ...)
```

**Arguments**

|                           |  |
|---------------------------|--|
| <code>partition</code>    | Partition (clustering) corresponds to the training set samples. Cluster labels in the partition will be used as covariate to construct SF. |
| <code>surv.time</code>    | A numeric vector contains the follow-up information of patients in <i>partition</i> , will be used as response to construct SF.            |
| <code>status</code>       | A binary vector contains survival status of patients in the partition, normally 0=alive, 1=dead, will be used as response to construct SF. |
| <code>te.partition</code> | Partition to be evaluated  |
| <code>te.surv.time</code> | A numeric vector contains the follow-up information of patient's in <i>te.partition</i>  |
| <code>te.status</code>    | A binary vector contains survival status of patients in <i>te.partition</i> , normally 0=alive, 1=dead.                                    |
| <code>...</code>          | Arguments for <code>rfsrc</code> from the <code>randomForestSRC</code> package   |

**Value**

A vector of error rates. Length is equal the number of trees constructed. Default is 1000.

**Author(s)**

Askar Obulkasim

**References**

Ishwaran,H. et al., (2008). "Random survival forest", *Ann. App. Statist.*, 2, 841-860.  
 Obulkasim,A. et al., (2013). "Semi-supervised adaptive-height snipping of the Hierarchical Clustering tree", submitted.

**See Also**

[rfsrc](#)

**Examples**

```

data(BullingerLeukemia)
em <- BullingerLeukemia[["em"]]
st <- BullingerLeukemia[["surv.time"]]
event <- BullingerLeukemia[["status"]]
cl <- HCsnipper(em[, 1:60], min = 5)
pred <- cluster_pred(X = em[, 1:100], partition = cl$partitions[1, ], surv.time = st[1:60],
                    status = event[1:60], te.index = 61:100)
Err <- RSF_eval(cl$partitions[1, ], st[1:60], event[1:60], pred, st[61:100], event[61:100])

```

surv\_measure

*Cluster quality evaluation using follow-up data***Description**

Function to evaluate the overall quality of a given partition using follow-up data. A partition (clustering) is composed of non-overlapping clusters.

**Usage**

```
surv_measure(parti, surv.time, status, method = "BIC")
```

**Arguments**

|           |  |
|-----------|--|
| parti     | A partition to be evaluated.   |
| surv.time | A numeric vector contains follow-up time of patient's in $x$   |
| status    | A binary vector contains survival status of patients in $x$ , normally 0=alive, 1=dead.  |
| method    | Type of partition evaluation measures to use for assessing the relationship between follow-up and a partition. Default is <i>BIC</i> . |

**Details**

This function fits a Cox model using follow-up data as response and cluster labels in the partition as covariate. The likelihood from the fitted model further used to calculate the modified *AIC* or *BIC*. See references for more details. Note that, for convenience in later usage, returned value is multiplied by -1 inside the function so that large value denotes good quality partition.

**Value**

A numeric value representing the quality of partition under consideration in terms of follow-up.

**Author(s)**

Askar Obulkasim

## References

- Liang,H. and Zou,G.H. (2008). "Improved AIC selection strategy for survival analysis", *Comput Stat Anal.*, 52, 2538-2548.
- Volinsky,T.C. and Raftery,A.E. (2000). "Baysian information criteria for censored survival models", *Biometrics*, 56, 256-262.
- Obulkasim,A. et al., (2013). "Semi-supervised adaptive-height snipping of the Hierarchical Clustering tree", submitted.

## See Also

[measure](#)

## Examples

```
data(BullingerLeukemia)
em <- BullingerLeukemia[["em"]]
st <- BullingerLeukemia[["surv.time"]]
event <- BullingerLeukemia[["status"]]
cl <- HCsnipper(em[, 1:60], min = 5)
result <- apply(cl[["partitions"]], 1, function(x) surv_measure(x, st[1:60], event[1:60]))
```

---

TcgaGBM

*Glioblastoma multiforme gene expression data*

---

## Description

The subset of latest version of TCGA glioblastoma (GBM) level 3 gene expression data with partial clinical info. Contains expression profiles of 120 samples measured on 3000 genes. Clinical data includes follow-up and type of drugs patients have been administered.

## Usage

```
data(TcgaGBM)
```

## Format

An object of list class

## Note

The complete dataset is available at the TCGA data portal <https://tcga-data.nci.nih.gov/tcga/>.

## Source

The Cancer Genome Atlas Network (2008), "Comprehensive genomic characterisation defines human glioblastoma genes and core pathways", *Nature*, 490, 61-70.

**Examples**

```
data(TcgaGBM)
names(TcgaGBM)
```

---

|              |  |
|--------------|--|
| TwoHC_assign | <i>Function to assign new samples to one of the two given hierarchical clustering trees in a semi-supervised way</i> |
|--------------|--|

---

**Description**

For given molecular data sets from two non-overlapping groups of patients, this functions constructs two independent HC trees and assigns new samples to one of them in semi-supervised way. See details.

**Usage**

```
TwoHC_assign(X, index1, index2, new.X, dis.method = "cor", link.method = "ward",
             minclus = 4, maxmiss = 30, surv.time, status, method1 = "BIC",
             method2 = "g2")
```

**Arguments**

|             |  |
|-------------|--|
| X           | An object of class <a href="#">ExpressionSet</a> or data matrix from which two HC trees to be derived. Columns are assumed to represent the samples, and rows represent the sample's features. Missing values are allowed. |
| index1      | Column indices of patients in X correspond to the first group.   |
| index2      | Column indices of patients in X correspond to the second group.  |
| new.X       | An object of class <a href="#">ExpressionSet</a> or data matrix corresponds to new samples. Columns are assumed to represent the samples, and rows represents the sample's features. Missing values are allowed.           |
| dis.method  | The distance measure to be used. This must be one of method acceptable for <a href="#">dist</a> function or the Pearson correlation (default).   |
| link.method | The agglomeration method to be used. This should be one of "ward" (default), "single", "complete", "average", "mcquitty", "median" or "centroid".  |
| minclus     | The minimum number of samples allowed to form a cluster. This parameter inversely proportional to the number of partition returned from a HC tree. e.g. a large value returns small number of partitions, and vice versa.  |
| maxmiss     | Maximum percentage of missing values per row in X.   |
| surv.time   | A numeric vector contains follow-up information of patient's in X  |
| status      | A binary vector contains survival status of patients in X, normally 0=alive, 1=dead.   |
| method1     | Type of partition evaluation measures to use for assessing the relationship between follow-up and a partition. Default is "BIC".   |
| method2     | Type of Partition evaluation measure to use for assessing the relationship between data matrix X and a partition. Default is Goodman and Kruskal index "g2".   |

## Details

Say molecular profiles of two groups patients (without overlap) treated with two different drugs or the same drugs in different combinations are available. Besides that, their follow-up information are also given. When a new patient comes in (for which only molecular profiles are available), question will be to which group this patient should be assigned so that he/she will benefit most by the type of treatment this group received.

This function is designed for this problem. it works as follows: first, two independent HC trees will be derived from given data; second, partitions are extracted and the optimal partition is selected from each HC tree, separately; third, new patient's molecular profile is compared with each cluster in each optimal partition to calculate average similarity and identify two most similar clusters (competing clusters) from the two HC trees; finally, new sample is assigned to one of the two competing clusters which has better overall survival.

## Value

A list object contains following components:

|                |  |
|----------------|--|
| hc1            | HC tree derived from the data corresponds to the first group.  |
| hc2            | HC tree derived from the data corresponds to the second group.   |
| partitions.hc1 | A matrix includes partitions extracted from <i>hc1</i> . Rows represent partitions and columns represent samples.  |
| partitions.hc2 | A matrix includes partitions extracted from <i>hc2</i> . Rows represent partitions and columns represent samples.  |
| best.hc1       | Optimal partition found on the <i>hc1</i>  |
| best.hc2       | Optimal partition found on the <i>hc2</i>  |
| score.hc1      | A matrix with two columns. The first column contains the quality scores of <i>partitions.hc1</i> calculated using the follow-up data. The second column contains the quality scores of <i>partition.hc1</i> calculated by using <i>X</i> .   |
| score.hc2      | The same as <i>score.hc1</i> , but for <i>partitions.hc2</i> .   |
| Assign         | A matrix with three columns. The first column contains the indices of HC trees to which a test sample was assigned. The second column contains the indices of clusters in <i>best.hc1</i> to which a test sample was most similar. The third column contains the indices of clusters in <i>best.hc2</i> to which a test sample was most similar. |
| surv.time      | The same as input  |
| status         | The same as input  |
| index1         | The same as input  |
| index2         | The same as input  |
| new.X          | The same as input  |
| X              | The same as input  |
| method1        | The same as input  |
| method2        | The same as input  |
| minclus        | The same as input  |

**Author(s)**

Askar Obulkasim

**References**

- Harrel,E.F. et al., (1982). "Evaluating the yield of medical tests", *JAMA*, 247, 2543-2546.
- Obulkasim,A. et al., (2011). "Stepwise classification of cancer samples using clinical and molecular data", *BMC Bioinformatics*, 12, 422.
- Troyanskaya,O. et al., (2001). "Missing value estimation methods for DNA microarrays". *Bioinformatics*, 17, 520-525.
- Obulkasim,A. et al., (2013). "Semi-supervised adaptive-height snipping of the Hierarchical Clustering tree", submitted.

**See Also**

See also [TwoHC\\_perm](#), [cluster\\_pred](#)

**Examples**

```
data(TcgaGBM)
id1 <- which(TcgaGBM[["drugs"]] == "Avastin")
id2 <- which(TcgaGBM[["drugs"]] == "Temodar")
em <- TcgaGBM[["em"]]
st <- TcgaGBM[["surv.time"]]
event <- TcgaGBM[["status"]]
result <- TwoHC_assign(X = em[,c(id1[1:50], id2[1:50])], index1 = 1:50, index2 = 51:100,
                      new.X = em[, c(id1[51:60], id2[51:60])], minclus = 4,
                      surv.time = st[c(id1[1:50], id2[1:50])],
                      status = event[c(id1[1:50], id2[1:50])])
```

---

|            |   |
|------------|---|
| TwoHC_perm | <i>Function to assess the significance of group assignment from</i><br><a href="#">TwoHC_assign</a> |
|------------|---|

---

**Description**

Function to evaluate the significance of the group assignments generated by [TwoHC\\_assign](#).

**Usage**

```
TwoHC_perm(TwoHC, nperm = 1000)
```

**Arguments**

|       |   |
|-------|---|
| TwoHC | Output from the 'TwoHC_assign' function |
| nperm | The number of permutations.             |

## Details

Significance of group assignment for each patient is calculated as follows: for a given patient, examine the previously found optimal partition in each HC tree and identify two clusters to which this patient is most similar. Say these two competing clusters are *cluster1* and *cluster2* of size *n1* and *n2*, respectively. Create a binary vector (call it *x*) of size *n1* + *n2* which has *n2* ones and *n1* zeros. Construct a Cox model using follow-up information of samples in *cluster1* and *cluster2* as response, and *x* as covariate. The absolute value of the estimated group parameter ('beta\_obs') in the Cox model that compares the survival times of the other patients in the two competing clusters expresses the predicted gain in survival from the assignment by 'TwoHC\_assign' with respect to random. The *beta\_obs* will be transformed to

$$r_{obs}^i = \exp\left(-|\hat{\beta}_{obs}^i|\right),$$

which quantifies the gain in relative risk in the Cox model. The problem is that this is biased, because 'TwoHC\_assign' already used 'beta\_obs'. Hence, even when the two groups would be equally good for the molecular profiles in the two competing clusters, we obtain 'r\_obs' < 1. To correct for this bias this function uses a permutation argument. For each new patient it applies 'nperm' permutations of the survival data among the two competing clusters. As above we compute 'r\_perm' for each permutation which contains the same bias as 'r\_obs'. Let  $Z_i = \text{median}(r\_perm(1), \dots, r\_perm(nperm))$ , then risk-ratio  $rr\_obs(i) = r\_obs(i) / Z_i$  quantifies the biased-corrected reduction in relative risk. The permuted version 'rr\_perm(i)' is defined analogously (see vignette). Finally it defines a test statistic:

$$T_{obs} = \frac{\frac{1}{n} \sum_{i=1}^n \log(rr_{obs}^i)}{stdev(\log(rr_{obs}^1, \dots, rr_{obs}^n))}$$

'T\_obs' compared with the background of its null-distribution as obtained by permutation to calculate *p*-value.

## Value

A list object contains following objects:

|            |  |
|------------|--|
| Obs.betas  | A numeric vector contains the coefficient from the Cox model corresponding to each test sample.  |
| Perm.betas | A matrix contains the coefficient from the Cox model trained with permuted follow-up data. Columns represent test samples, rows represent permutations |
| Ranks      | A numeric vector contains the rank of each observed coefficient among the <i>nperm</i> coefficients generated by permutations.                         |
| RiskRatios | A numeric vector contains the ratio of relative risk for the test set.   |
| Pvalue     | <i>p</i> -value of the overall group assignment.   |

## Author(s)

Askar Obulkasim

## References

- Obulkasim,A. et al., (2013). "Semi-supervised adaptive-height snipping of the Hierarchical Clustering tree", submitted.
- Harrel,E.F. et al., (1982). "Evaluating the yield of medical tests", *JAMA*, 247, 2543-2546.
- Obulkasim,A. et al., (2011). "Stepwise classification of cancer samples using clinical and molecular data", *BMC Bioinformatics*, 12, 422.

## See Also

See also [TwoHC\\_assign](#), [cluster\\_pred](#)

## Examples

```
data(TcgaGBM)
id1 <- which(TcgaGBM[["drugs"]] == "Avastin")
id2 <- which(TcgaGBM[["drugs"]] == "Temodar")
em <- TcgaGBM[["em"]]
st <- TcgaGBM[["surv.time"]]
event <- TcgaGBM[["status"]]
twoHC <- TwoHC_assign(X = em[,c(id1[1:50], id2[1:50])], index1 = 1:50, index2 = 51:100,
                      new.X = em[, c(id1[51:60], id2[51:60])], minclus = 4,
                      surv.time = st[c(id1[1:50], id2[1:50])],
                      status = event[c(id1[1:50], id2[1:50])])
result <- TwoHC_perm(twoHC, nperm = 100)

## Not run:
### Examples with a larger number of permutations (not run).
result <- TwoHC_perm(twoHC, nperm = 10000)
par(mfrow = c(1, 2))
plot(density(result[["Ranks"]]), xlab = "Ranks")
plot(density(result[["RiskRatios"]]), xlab = "Observed relative risk-ratios")

## End(Not run)
```

# Index

BullingerLeukemia, [2](#)

cluster.stats, [8](#), [9](#)

cluster\_pred, [3](#), [17](#), [19](#)

dist, [3](#), [7](#), [15](#)

EnvioPlot, [5](#)

ExpressionSet, [3](#), [5](#), [7](#), [15](#)

hclust, [7](#)

HCsnip-package, [2](#)

HCsnipper, [6](#)

hdEntropy, [6](#)

impute.knn, [4](#), [7](#)

measure, [8](#), [14](#)

perm\_test, [9](#)

rfsrc, [12](#)

RSF\_eval, [12](#)

surv\_measure, [9](#), [13](#)

TcgaGBM, [14](#)

TwoHC\_assign, [5](#), [15](#), [17](#), [19](#)

TwoHC\_perm, [11](#), [17](#), [17](#)