

Package: RHclust (via r-universe)

September 3, 2024

Type Package

Title Vector in Partition

Version 2.0.0

Maintainer Joseph Handwerker <jkhndwrk@memphis.edu>

Description Non-parametric clustering of joint pattern multi-genetic/epigenetic factors. This package contains functions designed to cluster subjects based on gene features including single nucleotide polymorphisms (SNPs), DNA methylation (CPG), gene expression (GE), and covariate data. The novel concept follows the general K-means (Hartigan and Wong (1979) <doi:10.2307/2346830> framework but uses weighted Euclidean distances across the gene features to cluster subjects. This approach is unique in that it attempts to capture all pairwise interactions in an effort to cluster based on their complex biological interactions.

License Artistic-2.0

Encoding UTF-8

Imports Runuran,graphics,stats,utils

NeedsCompilation no

Author Joseph Handwerker [aut, cre], Lauren Sobral [ctb], Meredith Ray [aut, ctb]

Repository CRAN

Date/Publication 2023-08-15 07:20:21 UTC

Contents

BinaryClass	2
SimData	3
VIP	5
VIPcov	8
VIPnoCPG	11
VIPnoSNP	13

BinaryClass	<i>Binary Classification</i>
-------------	------------------------------

Description

A confusion matrix but allows for analysis of non-equal level data classifications.

Usage

```
BinaryClass(x)
```

Arguments

x Can be a data frame dimensions at least 2 rows and 2 columns meant to represent observed and predicted values where the observed (true) values are in the first column and predicted columns in the second column.

Details

BinaryClass() is similar to a confusion matrix with binary classification outputs. The true positive values per column are identified based on the maximum number of assignments per category.

Value

Table	the results of 'table()' on 'x'
Accuracy	overall accuracy of classification
CI	confidence interval of overall accuracy using Clopper-Pearson Interval
Group Measures	the sensitivity, specificity, positive predictive value, negative predictive value, prevalence detection rate, detection prevalence, and balanced accuracy for each class

Author(s)

jkhdwrk@memphis.edu

Examples

```
# Basic example
true = c(rep(1,5), rep(2,5), rep(3,5), rep(4,5))
pred = c(rep(1,4), 4, rep(2,5), 2, rep(3,4), 1, rep(4,4))
df = cbind(true, pred)
BinaryClass(df)

true = c(rep(1,5), rep(2,5), rep(3,5), rep(4,5))
pred = c(rep(1,5), rep(2,5), rep(3,10))
df = cbind(true, pred)
```

```

BinaryClass(df)

sd = SimData(k = c(10,40,50))
out = VIP(sd, v = 3, optimize = 'elbow', nstart = 5)
df = out$`BC Test`
BinaryClass(df)

## Looping through different clusters

sd = SimData(seed = 1, gene = 1)
acc = NULL
for (i in 1:5){
  out = VIP(sd, v = i, optimize = 'off', nstart = 5)
  acc[i] = BinaryClass(out$`BC Test`)$Accuracy
}

plot(acc, type = 'b', main = 'Accuracy Comparison', xlab = 'Clusters', ylab = 'Acc')

```

SimData

GE, CPG, SNP, and Covariate Simulated Data

Description

Simulated data generator containing continuous variables representing gene expression (GE) data and DNA methylation data as M-values (GPG), and categorical variable representing single nucleotide polymorphisms (SNP). GE and CPG data are simulated from a normal distribution and SNP data is simulated from a multinomial distribution. Covariate data can have uniformly distributed data or normally distributed. Cluster separation can also be distinct or non-distinct.

Usage

```

SimData(seed = NULL, gene = 36,
        k = c(33,33,34),
        GEbar = 5, GESd = 0.5,
        CPGbar = 4, CPGsd = 0.5,
        SameCPG = FALSE, SameSNP = FALSE,
        ProbDist = NULL, SameGeneDist = TRUE,
        distinct=TRUE)

```

Arguments

seed	Set specified seed for reproducibility
gene	Numeric input that specifies the number genes

k	Cluster pattern/distribution across subjects formatted as a vector, i.e. <code>c(33,33,34)</code> representing 33 subjects in the first cluster, 33 in the second cluster, and 34 in the third cluster.
GEbar	Optional numeric input to change the mean distribution of GE data
GEsd	Optional numeric input to change the standard deviation of GE data
CPGbar	Optional numeric input to change the mean distribution of CPG data
CPGsd	Optional numeric input to change the standard deviation of CPG data
SameCPG	Logical value that if set to True sets the distribution of each CPG cluster around the same mean
SameSNP	Logical value that if set to True changes the probability distribution of SNPs to be the same per cluster
ProbDist	Optional list input that allows the change of SNP probability distributions per cluster. Default list stops at 5 cluster distributions. Default ProbDist = <code>list(c(0.50,0.25,0.25), c(0.20,0.55,0.25), c(0.30,0.15,0.55), c(0.20,0.50,0.30), c(0.45,0.20,0.35))</code>
SameGeneDist	Logical that set the covariates to follow the genetic clustering scheme where the data is uniformly distributed amongst each subgroup in each cluster. Or if set to FALSE sets the covariates to be noisy and have no discernible groups.
distinct	Logical that sets the covariate data clusters to be distinct at default, or non-distance when set to FALSE

Details

SimData creates simulated data that aims to represent real world data for gene expression (GE), DNA methylations (CPG), and single nucleotide polymorphisms (SNP). Covariate data is used to represent other potentially useful health data to further weight the genetics data. The goal of this function is to allow the user the ability to manipulate their data for testing of the `VIP()` or `VIPcov()` functions.

Value

Clusters	Vector of cluster assignment for each subject.
Vec	Numeric representation of values per cluster used for sensitivity measures.
GE	Simulated continuous data for GE. Means of each cluster changes by a factor of 5 with default standard deviation of 0.5.
CPG	Simulated continuous data for CPG. Means of each cluster changes by a factor of 4 with default standard deviation of 0.5.
SNP	Simulated categorical data for SNP.
GE_Index	Index names for GE.
CPG_Index	Index names for CPG.
SNP_Index	Index names for SNP.
Covariates	Simulated data for covariate data. By default contains 10 columns of categorical data and 10 columns of numeric data.

Author(s)

jkhndwrk@memphis.edu

Examples

```
# Generating simulated data
sd = SimData()

## Specifying seed, genes, and clusters
# sd = SimData(seed = 42, gene = 18, c(10,40,50))

# Specifying probability distribution of SNP data
l = list( c(10, 35, 55),
          c(60, 20, 20),
          c(25, 50, 25))
sd = SimData(1, g = 36, ProbDist = l)

sd = SimData(1, g = 36, c(33,33,34))

# Same SNP distribution across 2 clusters
ssnp = SimData(g = 36, k = c(10,90), SameSNP = TRUE)

# Same CpG distribution across 2 cluster
scpg = SimData(g = 36, k = c(10,90), SameCPG = TRUE)
```

 VIP

Vector in Partition

Description

Clustering of subjects based on similar patterns of gene expression, DNA methylation, and SNPs.

Usage

```
VIP(Simulated = NULL, SNP = NULL, CPG = NULL, GE = NULL,
     SNPname = NULL, CPGname = NULL, Gname = NULL, v,
     optimize = c('off', 'min', 'slope', 'elbow'),
     iter_max = 1000, nstart = 5, fit = c('aic', 'bic'),
     seed = NULL, type = c('Default', 'NoCPG', 'NoSNP'),
     ct = c('mean', 'median'), verbose = FALSE)
```

Arguments

Simulated	set to name of simulated data built from SimData(), else set to NULL for real data.
SNP	Data frame or data matrix containing categorical SNP data. Input must be in form of N x M, with N rows of subjects and M columns of SNPs. Rownames are permitted. Run SimData()\$SNP for examples.

CPG	Data frame or data matrix containing numeric CPG data. Input must be in form of $N \times M$, with N rows of subjects and M columns of CPG. Rownames are permitted. Run <code>SimData()\$CPG</code> for examples.
GE	Data frame or data matrix containing numeric GE data. Input must be in form of $N \times M$, with N rows of subjects and M columns of GE. Rownames are permitted. Run <code>SimData()\$GE</code> for examples.
SNPname	Names for SNP data. Data must be a data frame of $N \times 2$ dimensions with SNP sites as column 1, and GE indexes in column 2. Order of SNPs must match the order of the SNP columns in the argument SNP. See <code>SimData()\$SNP_Index</code> for examples.
CPGname	Names for CPG data. Data must be a data frame of $N \times 2$ dimensions with CPG sites as column 1, and GE indexes in column 2. Order of CPGs must match the order of the CPG columns in the argument GE. See <code>SimData()\$CPG_Index</code> for examples.
GEname	Names for GE data. Data must be a data frame of $N \times 2$ dimensions with GE sites as column 1, and GE indexes in column 2. Order of GEs must match the order of the GE columns in the argument GE. See <code>SimData()\$GE_Index</code> for examples.
v	Numeric scalar or vector of number for clusters, or a range of clusters with format <code>c(l,u)</code> for cluster $l:u$
optimize	Returned the optimal number of clusters. Input 'min' returns cluster assignment with lowest WSS for clusters in v. Input 'slope' indicates whether the algorithm should pick the lowest WSS value based on the first increasing slope. Input 'elbow' fits a line between the first and last fitted WSS and finds the corresponding cluster with the maximum distance to that line. All but 'slope' return plots.
iter_max	Maximum number of iterations allowed.
nstart	If <code>nstart > 1</code> , repetitive computations with random initializations are computed and the result with minimum <code>tot_dist</code> is returned.
fit	Penalizing factor for WSS of clusters. Can be set to either 'aic' or 'bic'.
seed	Optional input to sample the same initial cluster centers.
type	Optional input for special cases for data without CPGs or SNP inputs. Options include "Default", "NoSNP", or "NoCPG"
ct	Central tendency option for cluster assignment. Options include 'mean' or 'median'.
verbose	Logical whether information about the cluster procedure should be given.

Details

Similar to k-means and k-proto clustering, this algorithm computes clusters based on weighted factors of mixed data relative to genetic/epigenetic data. Clusters are assigned using summed euclidean distance of numerics (*GE* and *CPG*) weighted by matching categorical (*SNP*) data. Central tendency of numeric data can be set to either mean or median with input *ct*.

Data must be ordered such that rows in each data set correspond to the same subject and order of the indexes match the order of the columns in the data. The current algorithm does not allow for any missing data. The aim is for *GE*, *CPG*, and *SNP* data to be clustered into v groups such that

within sum of squares is minimized. If groups of clusters are close, the algorithm may not converge correctly and signals a warning if cluster size is reduced.

Optimization functionality was used for simulated data analysis, but is allowed for user exploratory analysis as well. *'min'* simply returns the lowest fitted WSS *fit* parameter. *'slope'* loops through clusters in v and returns the cluster based on the first increasing slope of fitted WSS. For example, if AIC output is $c(100,80,35,50)$, cluster 3 would be returned since the slope increases from 3 to 4. If there is no increasing slope, the *'min'* optimizer will be returned. *'elbow'* seeks to find the elbow of the plot based on saturation point. This worked the best for simulation studies but requires more clusters to make proper predictions, in our case it required a range of at least 5 clusters $c(1,5)$ to search to correctly identify the 3 simulated clusters. . For ease of exploratory analysis, $v=1$ is allowed.

Value

size	Number of subjects assigned to each cluster.
cluster	Vector of cluster assignment.
GECenters	Matrix of cluster centers for GE.
CPGCenters	Matrix of cluster centers for CPG.
SNPCenters	Matrix of cluster centers for SNP.
within	Vector of within cluster sum of squares with one component per cluster.
tot_within	Summed total of within-cluster sum of squares.
Moved	Number of iterations before convergence.
AIC	Value of tot_within with aic penalizer.
BIC	Value of tot_within with bic penalizer.
outputPlot	Returns the tot_within, aic, bic, and v values for plotting.

Author(s)

jkhndwrk@memphis.edu

References

Hartigan, J. A. and Wong, M. A. (1979). Algorithm AS 136: A K-means clustering algorithm. Applied Statistics, 28, 100–108. 10.2307/2346830.

Examples

```
## simple output of 3 clusters assignments
sd = SimData(1, g = 36, c(33,33,34))
VIPout = VIP(sd, v = 3)

# loop through clusters 1-10 and outputs plot of WSS, AIC, and BIC
VIPout = VIP(sd, v = c(1,10))

# loop through clusters 1-10 but picks first instance of increasing slope
VIPout = VIP(sd, v = c(1,10), optimize = 'slope')
```

```
# Individual inputs
sd = SimData(1, g = 36, k = c(33,33,34))
VIPout = VIP(SNP = sd$SNP, CPG = sd$CPG, GE = sd$GE,
             SNPname = sd$SNP_Index, CPGname = sd$CPG_Index,
             Gname = sd$GE_Index,
             v = c(1,5), optimize = 'off', nstart = 5)

## Varying clusters
sd = SimData(k = c(10,40,50))
out = VIP(sd, v = c(1,6), optimize = 'elbow', nstart = 30)
```

VIPcov

Vector in Partition with covariates

Description

Clustering of subjects based on similar patterns of gene expression, DNA methylation, and SNPs weighted by covariates.

Usage

```
VIPcov(Simulated = NULL, SNP = NULL, CPG = NULL, GE = NULL,
       SNPname = NULL, CPGname = NULL, Gname = NULL,
       covariates = NULL, v, lambda = NULL, gamma = 1,
       optimize = c('off', 'min', 'slope', 'elbow'),
       iter_max = 1000, nstart = 5, fit = c('aic', 'bic'),
       seed = NULL, type = c('Default', 'NoCPG', 'NoSNP'),
       ct = c('mean', 'median'), verbose = FALSE)
```

Arguments

Simulated	set to name of simulated data built from SimData(), else set to NULL for real data.
SNP	Data frame or data matrix containing categorical SNP data. Input must be in form of N x M, with N rows of subjects and M columns of SNPs. Rownames are permitted. Run SimData()\$SNP for examples.
CPG	Data frame or data matrix containing numeric CPG data. Input must be in form of N x M, with N rows of subjects and M columns of CPG. Rownames are permitted. Run SimData()\$CPG for examples.
GE	Data frame or data matrix containing numeric GE data. Input must be in form of N x M, with N rows of subjects and M columns of GE. Rownames are permitted. Run SimData()\$GE for examples.

SNPname	Names for SNP data. Data must be a data frame of Nx2 dimensions with SNP sites as column 1, and GE indexes in column 2. Order of SNPs must match the order of the SNP columns in the argument SNP. See SimData()\$SNP_Index for examples.
CPGname	Names for CPG data. Data must be a data frame of Nx2 dimensions with CPG sites as column 1, and GE indexes in column 2. Order of CPGs must match the order of the CPG columns in the argument GE. See SimData()\$CPG_Index for examples.
GEname	Names for GE data. Data must be a data frame of Nx2 dimensions with GE sites as column 1, and GE indexes in column 2. Order of GEs must match the order of the GE columns in the argument GE. See SimData()\$GE_Index for examples.
covariates	Data frame or matrix containing numeric and or categorical variables used for weighting the distance matrix of the genetics data.
v	Numeric scalar or vector of number for clusters, or a range of clusters with format c(l,u) for cluster l:u
lambda	Weighting factor for the covariate data. The default value is calculated by taking the average variance of the numeric variables divided by the average concentration of the factors for the categorical variables. When non-mixed type data are used, the value is set to 1 unless specified otherwise.
gamma	Weighting factor for the genetics distance measure. Default is set to 1
optimize	Returned the optimal number of clusters. Input 'min' returns cluster assignment with lowest WSS for clusters in v. Input 'slope' indicates whether the algorithm should pick the lowest WSS value based on the first increasing slope. Input 'elbow' fits a line between the first and last fitted WSS and finds the corresponding cluster with the maximum distance to that line. All but 'slope' return plots.
iter_max	Maximum number of iterations allowed.
nstart	If nstart > 1, repetitive computations with random initializations are computed and the result with minimum tot_dist is returned.
fit	Penalizing factor for WSS of clusters. Can be set to either 'aic' or 'bic'.
seed	Optional input to sample the same initial cluster centers.
type	Optional input for special cases for data without CPGs or SNP inputs. Options include "Default", "NoSNP", or "NoCPG"
ct	Central tendency option for cluster assignment. Options include 'mean' or 'median'.
verbose	Logical whether information about the cluster procedure should be given.

Details

This method is an extension of the VIP() function and allows for the clustering method to be further weighted by the covariate distance. The covariate distance measure is similar to that seen in *kproto()* where the euclidian measure of the numerics data is weighted by matching categorical variables and a weighting factor. Clusters are assigned based on these two distance measures where the Central tendency of numeric data can be set to either mean or median with input *ct*.

Data must be ordered such that rows in each data set correspond to the same subject and order of the indexes match the order of the columns in the data. The current algorithm does not allow for

any missing data. The aim is for *GE*, *CPG*, and *SNP* data to be clustered into v groups after being weighted by covariates such that within sum of squares is minimized. If groups of clusters are close, the algorithm may not converge correctly and signals a warning if cluster size is reduced.

Optimization functionality was used for simulated data analysis, but is allowed for user exploratory analysis as well. *'min'* simply returns the lowest fitted WSS *fit* parameter. *'slope'* loops through clusters in v and returns the cluster based on the first increasing slope of fitted WSS. For example, if AIC output is $c(100,80,35,50)$, cluster 3 would be returned since the slope increases from 3 to 4. If there is no increasing slope, the *'min'* optimizer will be returned. *'elbow'* seeks to find the elbow of the plot based on saturation point. This worked the best for simulation studies but requires more clusters to make proper predictions, in our case it required a range of at least 5 clusters $c(1,5)$ to search to correctly identify the 3 simulated clusters. For ease of exploratory analysis, $v=1$ is allowed. *'lambda'* is used as a weighting factor for the covariate distance measure, while *'gamma'* is the weighting factor for the genetics distance.

Value

size	Number of subjects assigned to each cluster.
cluster	Vector of cluster assignment.
GECenters	Matrix of cluster centers for GE.
CPGCenters	Matrix of cluster centers for CPG.
SNPCenters	Matrix of cluster centers for SNP.
Adjusted	Logical if the genetics distance measure is weighted by covariates.
Lambda	Lambda value used to weight the covariate distance measure.
within	Vector of within cluster sum of squares with one component per cluster.
tot_within	Summed total of within-cluster sum of squares.
Moved	Number of iterations before convergence.
AIC	Value of tot_within with aic penalizer.
BIC	Value of tot_within with bic penalizer.
outputPlot	Returns the tot_within, aic, bic, and v values for plotting.

Author(s)

jkhndwrk@memphis.edu

References

- Hartigan, J. A. and Wong, M. A. (1979). Algorithm AS 136: A K-means clustering algorithm. *Applied Statistics*, 28, 100–108. [10.2307/2346830](https://doi.org/10.2307/2346830).
- Z. Huang. (1998) Extensions to the k-means algorithm for clustering large data sets with categorical variables. *Data Mining and Knowledge Discovery*, 2:283–304, doi: [10.1023/A:1009769707641](https://doi.org/10.1023/A:1009769707641).
- Szpannek G (2018). “clustMixType: User-Friendly Clustering of Mixed-Type Data in R.” *The R Journal*, 200-208. doi:[10.32614/RJ-2018-048](https://doi.org/10.32614/RJ-2018-048), <https://doi.org/10.32614/RJ-2018-048>.

Examples

```
## simple output of 3 clusters assignments
sd = SimData(1, g = 36, c(33,33,34))
VIPout = VIPcov(sd, v = 3, covariates = sd$Covariates)

# loop through clusters 1-10 and outputs plot of WSS, AIC, and BIC
VIPout = VIPcov(sd, v = c(1,10), covariates = sd$Covariates)

# loop through clusters 1-10 but picks first instance of increasing slope
VIPout = VIPcov(sd, v = c(1,10), optimize = 'slope', covariates = sd$Covariates)

# Individual inputs
sd = SimData(1, g = 36, k = c(33,33,34))
VIPout = VIPcov(SNP = sd$SNP, CPG = sd$CPG, GE = sd$GE,
                SNPname = sd$SNP_Index, CPGname = sd$CPG_Index,
                GEname = sd$GE_Index,
                v = c(1,5), optimize = 'off', nstart = 5, covariates = sd$Covariates)

## Varying clusters
sd = SimData(k = c(10,40,50))
out = VIPcov(sd, v = c(1,6), optimize = 'elbow', nstart = 30, covariates = sd$Covariates)
```

VIPnoCPG

Vector in Partition without CPG data

Description

Clustering of subjects based on similar patterns of gene expression and SNPs.

Usage

```
VIPnoCPG(Simulated = NULL, SNP = NULL, GE = NULL,
          SNPname = NULL, GEname = NULL, v,
          optimize = c('off', 'min', 'slope', 'elbow'),
          iter_max = 1000, nstart = 5, fit = c('aic', 'bic'),
          seed = NULL, ct = c('mean', 'median'), verbose = FALSE)
```

Arguments

Simulated	set to name of simulated data built from SimData(), else set to NULL for real data.
SNP	Data frame or data matrix containing categorical SNP data. Input must be in form of N x M, with N rows of subjects and M columns of SNPs. Rownames are permitted. Run SimData()\$SNP for examples.

GE	Data frame or data matrix containing numeric GE data. Input must be in form of $N \times M$, with N rows of subjects and M columns of GE. Rownames are permitted. Run <code>SimData()\$GE</code> for examples.
SNPname	Names for SNP data. Data must be a data frame of $N \times 2$ dimensions with SNP sites as column 1, and GE indexes in column 2. Order of SNPs must match the order of the SNP columns in the argument SNP. See <code>SimData()\$SNP_Index</code> for examples.
GEname	Names for GE data. Data must be a data frame of $N \times 2$ dimensions with GE sites as column 1, and GE indexes in column 2. Order of GEs must match the order of the GE columns in the argument GE. See <code>SimData()\$GE_Index</code> for examples.
v	Numeric scalar or vector of number for clusters, or a range of clusters with format <code>c(l,u)</code> for cluster <code>l:u</code>
optimize	Returned the optimal number of clusters. Input 'min' returns cluster assignment with lowest WSS for clusters in v. Input 'slope' indicates whether the algorithm should pick the lowest WSS value based on the first increasing slope. Input 'elbow' fits a line between the first and last fitted WSS and finds the corresponding cluster with the maximum distance to that line. All but 'slope' return plots.
iter_max	Maximum number of iterations allowed.
nstart	If <code>nstart > 1</code> , repetitive computations with random initializations are computed and the result with minimum <code>tot_dist</code> is returned.
fit	Penalizing factor for WSS of clusters. Can be set to either 'aic' or 'bic'.
seed	Optional input to sample the same initial cluster centers.
ct	Central tendency option for cluster assignment. Options include 'mean' or 'median'.
verbose	Logical whether information about the cluster procedure should be given.

Details

The details are outlined in the main `VIP()` function. The only difference in this function is the absence of CPG data.

Value

size	Number of subjects assigned to each cluster.
cluster	Vector of cluster assignment.
GECenters	Matrix of cluster centers for GE.
SNPCenters	Matrix of cluster centers for SNP.
within	Vector of within cluster sum of squares with one component per cluster.
tot_within	Summed total of within-cluster sum of squares.
Moved	Number of iterations before convergence.
AIC	Value of <code>tot_within</code> with aic penalizer.
BIC	Value of <code>tot_within</code> with bic penalizer.
outputPlot	Returns the <code>tot_within</code> , <code>aic</code> , <code>bic</code> , and <code>v</code> values for plotting.

Author(s)

jkhndwrk@memphis.edu

References

Hartigan, J. A. and Wong, M. A. (1979). Algorithm AS 136: A K-means clustering algorithm. Applied Statistics, 28, 100–108. 10.2307/2346830.

Examples

```
# No CPG data
sd = SimData()
noCPGout = VIP(sd, v = c(1,5), optimize = 'off', nstart = 30, type = 'NoCPG')

noCPGout = VIPnoCPG(sd, v = c(1,5), optimize = 'off', nstart = 30)
```

 VIPnoSNP

Vector in Partition without SNP data

Description

Clustering of subjects based on similar patterns of gene expression and DNA methylation.

Usage

```
VIPnoSNP(Simulated = NULL, CPG = NULL, GE = NULL,
          CPGname = NULL, GEname = NULL, v,
          optimize = c('off', 'min', 'slope', 'elbow'),
          iter_max = 1000, nstart = 5, fit = c('aic', 'bic'),
          seed = NULL, ct = c('mean', 'median'), verbose = FALSE)
```

Arguments

Simulated	set to name of simulated data built from SimData(), else set to NULL for real data.
CPG	Data frame or data matrix containing numeric CPG data. Input must be in form of N x M, with N rows of subjects and M columns of CPG. Rownames are permitted. Run SimData()\$CPG for examples.
GE	Data frame or data matrix containing numeric GE data. Input must be in form of N x M, with N rows of subjects and M columns of GE. Rownames are permitted. Run SimData()\$GE for examples.
CPGname	Names for CPG data. Data must be a data frame of Nx2 dimensions with CPG sites as column 1, and GE indexes in column 2. Order of CPGs must match the order of the CPG columns in the argument GE. See SimData()\$CPG_Index for examples.

GEname	Names for GE data. Data must be a data frame of Nx2 dimensions with GE sites as column 1, and GE indexes in column 2. Order of GEs must match the order of the GE columns in the argument GE. See SimData()\$GE_Index for examples.
v	Numeric scalar or vector of number for clusters, or a range of clusters with format c(l,u) for cluster l:u
optimize	Returned the optimal number of clusters. Input 'min' returns cluster assignment with lowest WSS for clusters in v. Input 'slope' indicates whether the algorithm should pick the lowest WSS value based on the first increasing slope. Input 'elbow' fits a line between the first and last fitted WSS and finds the corresponding cluster with the maximum distance to that line. All but 'slope' return plots.
iter_max	Maximum number of iterations allowed.
nstart	If nstart > 1, repetitive computations with random initializations are computed and the result with minimum tot_dist is returned.
fit	Penalizing factor for WSS of clusters. Can be set to either 'aic' or 'bic'.
seed	Optional input to sample the same initial cluster centers.
ct	Central tendency option for cluster assignment. Options include 'mean' or 'median'.
verbose	Logical whether information about the cluster procedure should be given.

Details

The details are outlined in the main VIP() function. The only difference in this function is the absence of SNP data.

Value

size	Number of subjects assigned to each cluster.
cluster	Vector of cluster assignment.
GECenters	Matrix of cluster centers for GE.
CPGCenters	Matrix of cluster centers for CPG.
within	Vector of within cluster sum of squares with one component per cluster.
tot_within	Summed total of within-cluster sum of squares.
Moved	Number of iterations before convergence.
AIC	Value of tot_within with aic penalizer.
BIC	Value of tot_within with bic penalizer.
outputPlot	Returns the tot_within, aic, bic, and v values for plotting.

Author(s)

jkhdwrk@memphis.edu

References

Hartigan, J. A. and Wong, M. A. (1979). Algorithm AS 136: A K-means clustering algorithm. Applied Statistics, 28, 100–108. 10.2307/2346830.

Examples

```
# No SNP data
sd = SimData()
noSNPout = VIP(sd, v = c(1,5), optimize = 'off', nstart = 30, type = 'NoSNP')

noSNPout = VIPnoSNP(sd, v = c(1,5), optimize = 'off', nstart = 30)
```

Index

BinaryClass, [2](#)

SimData, [3](#)

VIP, [5](#)

VIPcov, [8](#)

VIPnoCPG, [11](#)

VIPnoSNP, [13](#)