

# Package ‘HTScluster’

January 20, 2025

**Type** Package

**Title** Clustering High-Throughput Transcriptome Sequencing (HTS) Data

**Version** 2.0.11

**Date** 2023-09-04

**Author** Andrea Rau, Gilles Celeux, Marie-Laure Martin-Magniette, Cathy Maugis-Rabusseau

**Maintainer** Andrea Rau <andrea.rau@jouy.inra.fr>

**Depends** R (>= 2.10.0)

**Imports** edgeR, plotrix, capushe, grDevices, graphics, stats

**Suggests** HTSFilter, Biobase

**Description** A Poisson mixture model is implemented to cluster genes from high-throughput transcriptome sequencing (RNA-seq) data. Parameter estimation is performed using either the EM or CEM algorithm, and the slope heuristics are used for model selection (i.e., to choose the number of clusters).

**License** GPL (>= 3)

**LazyLoad** yes

**NeedsCompilation** no

**Repository** CRAN

**Date/Publication** 2023-09-05 08:00:02 UTC

## Contents

HTScluster-package	2
highDimensionARI	3
HTSclusterUsersGuide	4
Init	5
logLikePoisMix	8
plot.HTScluster	11
PoisMixClus	13
PoisMixMean	18
PoisMixSim	20
probaPost	21
summary.HTScluster	23

---

HTScluster-package      *Clustering high throughput sequencing (HTS) data*

---

### Description

A Poisson mixture model is implemented to cluster genes from high-throughput transcriptome sequencing (RNA-seq) data. Parameter estimation is performed using either the EM or CEM algorithm, and the slope heuristics are used for model selection (i.e., to choose the number of clusters).

### Author(s)

Andrea Rau, Gilles Celeux, Marie-Laure Martin-Magniette, Cathy Maugis-Rabusseau

Maintainer: Andrea Rau

### References

Rau, A., Maugis-Rabusseau, C., Martin-Magniette, M.-L., Celeux G. (2015). Co-expression analysis of high-throughput transcriptome sequencing data with Poisson mixture models. *Bioinformatics*, 31(9):1420-1427.

Rau, A., Celeux, G., Martin-Magniette, M.-L., Maugis-Rabusseau, C. (2011) Clustering high-throughput sequencing data with Poisson mixture models. Inria Research Report 7786. Available at <https://inria.hal.science/inria-00638082>.

### Examples

```
set.seed(12345)

## Simulate data as shown in Rau et al. (2011)
## Library size setting "A", high cluster separation
## n = 2000 observations

simulate <- PoisMixSim(n = 200, libsize = "A", separation = "high")
y <- simulate$y
conds <- simulate$conditions

## Run the PMM model for g = 3
## "TC" library size estimate, EM algorithm

run <- PoisMixClus(y, g=3, conds=conds, norm="TC")

## Estimates of pi and lambda for the selected model

pi.est <- run$pi
lambda.est <- run$lambda

## Not run: PMM for 4 total clusters, with one fixed class
```

```

## "TC" library size estimate, EM algorithm
##
## run <- PoisMixClus(y, g = 3, norm = "TC", conds = conds,
##   fixed.lambda = list(c(1,1,1)))
##
##
## Not run: PMM model for 4 clusters, with equal proportions
## "TC" library size estimate, EM algorithm
##
## run <- PoisMixClus(y, g = 4, norm = "TC", conds = conds,
##   equal.proportions = TRUE)
##
##
## Not run: PMM model for g = 1, ..., 10 clusters, Split Small-EM init
##
## run1.10 <- PoisMixClusWrapper(y, gmin = 1, gmax = 10, conds = conds,
## norm = "TC")
##
##
## Not run: PMM model for g = 1, ..., 10 clusters, Small-EM init
##
## run1.10bis <- <- PoisMixClusWrapper(y, gmin = 1, gmax = 10, conds = conds,
## norm = "TC", split.init = FALSE)
##
##
## Not run: previous model equivalent to the following
##
## for(K in 1:10) {
## run <- PoisMixClus(y, g = K, conds = conds, norm = "TC")
## }

```

---

highDimensionARI

*Calculate ARI for high-dimensional data via data splits*


---

### Description

This function is used to calculate Adjusted Rand Index (ARI) values for high-dimensional data.

### Usage

```
highDimensionARI(x, y, splits = 2, verbose = FALSE)
```

### Arguments

x	Vector of classification labels
y	Vector of classification labels
splits	Number of subsets data should be split into
verbose	TRUE if verbose output is desired

**Value**

Value of Adjusted Rand Index for samples  $x$  and  $y$

**Author(s)**

Andrea Rau

**References**

Rau, A., Maugis-Rabusseau, C., Martin-Magniette, M.-L., Celeux G. (2015). Co-expression analysis of high-throughput transcriptome sequencing data with Poisson mixture models. *Bioinformatics*, 31(9):1420-1427.

Rau, A., Celeux, G., Martin-Magniette, M.-L., Maugis-Rabusseau, C. (2011). Clustering high-throughput sequencing data with Poisson mixture models. Inria Research Report 7786. Available at <https://inria.hal.science/inria-00638082>.

---

HTSclusterUsersGuide *View HTScluster User's Guide*

---

**Description**

Finds the location of the HTScluster User's Guide and optionally opens it.

**Usage**

```
HTSclusterUsersGuide(view=TRUE)
```

**Arguments**

`view` logical, should the document be opened using the default PDF document reader?

**Details**

The function `vignette("HTScluster")` will find the short HTScluster Vignette which describes how to obtain the HTScluster User's Guide. The User's Guide is not itself a true vignette because it is not automatically generated using [Sweave](#) during the package build process. This means that it cannot be found using `vignette`, hence the need for this special function.

If the operating system is other than Windows, then the PDF viewer used is that given by `Sys.getenv("R_PDFVIEWER")`. The PDF viewer can be changed using `Sys.putenv(R_PDFVIEWER=)`.

Note that this function was adapted from that defined by Gordon Smyth in the `edgeR` package.

**Value**

Character string giving the file location. If `view=TRUE`, the PDF document reader is started and the User's Guide is opened, as a side effect.

**Author(s)**

Gordon Smyth

**See Also**[system](#)**Examples**

```
# To get the location:
HTSClusterUsersGuide(view=FALSE)
# To open in pdf viewer:
## Not run: HTSClusterUsersGuide()
```

---

**Init***Parameter initialization for a Poisson mixture model.*

---

**Description**

These functions implement a variety of initialization methods for the parameters of a Poisson mixture model: the Small EM initialization strategy (`emInit`) described in Rau et al. (2011), a K-means initialization strategy (`kmeanInit`) that is itself used to initialize the small EM strategy, the splitting small-EM initialization strategy (`splitEMInit`) based on that described in Papastamoulis et al. (2014), and a function to initialize a small-EM strategy using the posterior probabilities (`probaPostInit`) obtained from a previous run with one fewer cluster following the splitting strategy.

**Usage**

```
emInit(y, g, conds, norm, alg.type = "EM",
       init.runs, init.iter, fixed.lambda, equal.proportions, verbose)

kmeanInit(y, g, conds, norm, fixed.lambda,
          equal.proportions)

splitEMInit(y, g, conds, norm, alg.type, fixed.lambda,
            equal.proportions, prev.labels, prev.probaPost, init.runs,
            init.iter, verbose)

probaPostInit(y, g, conds, norm, alg.type = "EM",
              fixed.lambda, equal.proportions, probaPost.init, init.iter,
              verbose)
```

**Arguments**

<code>y</code>	$(n \times q)$ matrix of observed counts for $n$ observations and $q$ variables
<code>g</code>	Number of clusters. If <code>fixed.lambda</code> contains a list of lambda values to be fixed, $g$ corresponds to the number of clusters in addition to those fixed.
<code>conds</code>	Vector of length $q$ defining the condition (treatment group) for each variable (column) in <code>y</code>
<code>norm</code>	The type of estimator to be used to normalize for differences in library size: (“TC” for total count, “UQ” for upper quantile, “Med” for median, “DESeq” for the normalization method in the DESeq package, and “TMM” for the TMM normalization method (Robinson and Oshlack, 2010). Can also be a vector (of length $q$ ) containing pre-estimated library size estimates for each sample.
<code>alg.type</code>	Algorithm to be used for parameter estimation (“EM” or “CEM” for the EM or CEM algorithms, respectively)
<code>init.runs</code>	In the case of the Small-EM algorithm, the number of independent runs to be performed. In the case of the splitting Small-EM algorithm, the number of cluster splits to be performed in the splitting small-EM initialization.
<code>init.iter</code>	The number of iterations to run within each Small-EM algorithm
<code>fixed.lambda</code>	If one (or more) clusters with fixed values of lambda is desired, a list containing vectors of length $d$ (the number of conditions). Note that the values of lambda chosen must satisfy the constraint noted in the technical report.
<code>equal.proportions</code>	If TRUE, the cluster proportions are set to be equal for all clusters. Default is FALSE (unequal cluster proportions)
<code>prev.labels</code>	A vector of length $n$ of cluster labels obtained from the previous run ( $g-1$ clusters)
<code>prev.probaPost</code>	An $n \times (g-1)$ matrix of the conditional probabilities of each observation belonging to each of the $g-1$ clusters from the previous run
<code>probaPost.init</code>	An $n \times g$ matrix of the conditional probabilities of each observation belonging to each of the $g$ clusters following the splitting strategy in the <code>splitEMInit</code> function
<code>verbose</code>	If TRUE, include verbose output

**Details**

In practice, the user will not directly call the initialization functions described here; they are indirectly called for a single number of clusters through the `PoisMixClus` function (via `init.type`) or via the `PoisMixClusWrapper` function for a sequence of cluster numbers (via `gmin.init.type` and `split.init`).

To initialize parameter values for the EM and CEM algorithms, for the Small-EM strategy (Bieracki et al., 2003) we use the `emInit` function as follows. For a given number of independent runs (given by `init.runs`), the following procedure is used to obtain parameter values: first, a K-means algorithm (MacQueen, 1967) is run to partition the data into  $g$  clusters ( $\hat{z}^{(0)}$ ). Second, initial parameter values  $\pi^{(0)}$  and  $\lambda^{(0)}$  are calculated (see Rau et al. (2011) for details). Third, a given number of iterations of an EM algorithm are run (defined by `init.iter`), using  $\pi^{(0)}$  and  $\lambda^{(0)}$  as initial

values. Finally, among the `init.runs` sets of parameter values, we use  $\hat{\lambda}$  and  $\hat{\pi}$  corresponding to the highest log likelihood or completed log likelihood to initialize the subsequent full EM or CEM algorithms, respectively.

For the splitting small EM initialization strategy, we implement an approach similar to that described in Papastamoulis et al. (2014), where the cluster from the previous run (with  $g-1$  clusters) with the largest entropy is chosen to be split into two new clusters, followed by a small EM run as described above.

### Value

<code>pi.init</code>	Vector of length $g$ containing the estimate for $\hat{\pi}$ corresponding to the highest log likelihood (or completed log likelihood) from the chosen initialization strategy.
<code>lambda.init</code>	$(d \times g)$ matrix containing the estimate of $\hat{\lambda}$ corresponding to the highest log likelihood (or completed log likelihood) from the chosen initialization strategy, where $d$ is the number of conditions and $g$ is the number of clusters.
<code>lambda</code>	$(d \times g)$ matrix containing the estimate of $\hat{\lambda}$ arising from the splitting initialization and small EM run for a single split, where $d$ is the number of conditions and $g$ is the number of clusters.
<code>pi</code>	Vector of length $g$ containing the estimate for $\hat{\pi}$ arising from the splitting initialization and small EM run for a single split, where $g$ is the number of clusters.
<code>log.like</code>	Log likelihood arising from the splitting initialization and small EM run for a single split.

### Author(s)

Andrea Rau

### References

- Anders, S. and Huber, W. (2010) Differential expression analysis for sequence count data. *Genome Biology*, **11**(R106), 1-28.
- Biernacki, C., Celeux, G., Govaert, G. (2003) Choosing starting values for the EM algorithm for getting the highest likelihood in multivariate Gaussian mixture models. *Computational Statistics and Data Analysis*, **41**(1), 561-575.
- MacQueen, J. B. (1967) Some methods for classification and analysis of multivariate observations. In *Proceedings of the 5th Berkeley Symposium on Mathematical Statistics and Probability*, number 1, pages 281-297. Berkeley, University of California Press.
- Papastamoulis, P., Martin-Magniette, M.-L., and Maugis-Rabuseau, C. (2014). On the estimation of mixtures of Poisson regression models with large number of components. *Computational Statistics and Data Analysis*: 3rd special Issue on Advances in Mixture Models, DOI: 10.1016/j.csda.2014.07.005.
- Rau, A., Celeux, G., Martin-Magniette, M.-L., Maugis-Rabuseau, C. (2011). Clustering high-throughput sequencing data with Poisson mixture models. Inria Research Report 7786. Available at <https://inria.hal.science/inria-00638082>.
- Rau, A., Maugis-Rabuseau, C., Martin-Magniette, M.-L., Celeux G. (2015). Co-expression analysis of high-throughput transcriptome sequencing data with Poisson mixture models. *Bioinformatics*, **31**(9):1420-1427.

Robinson, M. D. and Oshlack, A. (2010) A scaling normalization method for differential expression analysis of RNA-seq data. *Genome Biology*, **11**(R25).

### See Also

[PoisMixClus](#) for Poisson mixture model estimation for a given number of clusters, [PoisMixClusWrapper](#) for Poisson mixture model estimation and model selection for a sequence of cluster numbers.

### Examples

```
set.seed(12345)

## Simulate data as shown in Rau et al. (2011)
## Library size setting "A", high cluster separation
## n = 500 observations

simulate <- PoisMixSim(n = 500, libsize = "A", separation = "high")
y <- simulate$y
conds <- simulate$conditions

## Calculate initial values for lambda and pi using the Small-EM
## initialization (4 classes, PMM-II model with "TC" library size)
##
## init.values <- emInit(y, g = 4, conds,
##   norm = "TC", alg.type = "EM",
##   init.runs = 50, init.iter = 10, fixed.lambda = NA,
##   equal.proportions = FALSE, verbose = FALSE)
## pi.init <- init.values$pi.init
## lambda.init <- init.values$lambda.init
```

---

logLikePoisMix

*Log likelihood calculation for a Poisson mixture model*

---

### Description

Functions to calculate the log likelihood for a Poisson mixture model, the difference in log likelihoods for two different sets of parameters of a Poisson mixture model or the log-likelihood for each observation.

### Usage

```
logLikePoisMix(y, mean, pi)
logLikePoisMixDiff(y, mean.new, pi.new, mean.old, pi.old)
myLogLikePoisMixObs(y, conds, s, lambda, pi)
```



**Arguments**

<code>y</code>	$(n \times q)$ matrix of observed counts for $n$ observations and $q$ variables
<code>mean</code>	List of length $g$ containing the $(n \times q)$ matrices of conditional mean expression for all observations, as calculated by the <code>PoisMixMean</code> function, where $g$ represents the number of clusters
<code>mean.new</code>	List of length $g$ containing the $(n \times q)$ matrices of conditional mean expression for all observations for one set of parameters, as calculated by the <code>PoisMixMean</code> function, where $g$ represents the number of clusters
<code>mean.old</code>	List of length $g$ containing the $(n \times q)$ matrices of conditional mean expression for all observations for another set of parameters, as calculated by the <code>PoisMixMean</code> function, where $g$ represents the number of clusters
<code>pi.new</code>	Vector of length $g$ containing one estimate for $\hat{\pi}$
<code>pi.old</code>	Vector of length $g$ containing another estimate for $\hat{\pi}$
<code>pi</code>	Vector of length $g$ containing estimate for $\hat{\pi}$
<code>conds</code>	Vector of length $q$ defining the condition (treatment group) for each variable (column) in <code>y</code>
<code>s</code>	Estimate of normalized per-variable library size
<code>lambda</code>	$(d \times g)$ matrix containing the current estimate of lambda, where $d$ is the number of conditions (treatment groups) and $g$ is the number of clusters

**Details**

The `logLikePoisMixDiff` function is used to calculate the difference in log likelihood for two different sets of parameters in a Poisson mixture model; it is used to determine convergence in the EM algorithm run by the `PoisMixClus` function. The `logLikePoisMix` function (taken largely from the `mylogLikePoisMix` function from the `poisson.glm.mix` R package) calculates the log likelihood for a given set of parameters in a Poisson mixture model and is used in the `PoisMixClus` function for the calculation of the BIC and ICL. The `mylogLikePoisMixObs` function calculates the log likelihood per observation for a given set of parameters in a Poisson mixture model.

**Value**

ll (Depending on the context), the log likelihood, difference in log likelihoods for two different sets of parameters, or per-observation log-likelihood

**Note**

In the `logLikePoisMixDiff` function, we make use of the alternative mass function for a Poisson density proposed by Loader (2000) to avoid computational difficulties. The `logLikePoisMixDiff` function returns a default value of 100 if one or both of the log likelihoods associated with the two parameter sets takes on a value of  $-\infty$ .

**Author(s)**

Andrea Rau

## References

- Loader, C. (2000) Fast and accurate computation of binomial probabilities. Available at <https://lists.gnu.org/archive/html/octave-maintainers/2011-09/pdfK0uKOST642.pdf>.
- Rau, A., Maugis-Rabusseau, C., Martin-Magniette, M.-L., Celeux G. (2015). Co-expression analysis of high-throughput transcriptome sequencing data with Poisson mixture models. *Bioinformatics*, 31(9):1420-1427.
- Rau, A., Celeux, G., Martin-Magniette, M.-L., Maugis-Rabusseau, C. (2011) Clustering high-throughput sequencing data with Poisson mixture models. Inria Research Report 7786. Available at <https://inria.hal.science/inria-00638082>.

## See Also

[PoisMixClus](#) for Poisson mixture model estimation and model selection; [PoisMixMean](#) to calculate the per-cluster conditional mean of each observation

## Examples

```
set.seed(12345)

## Simulate data as shown in Rau et al. (2011)
## Library size setting "A", low cluster separation
## n = 200 observations

simulate <- PoisMixSim(n = 200, libsize = "A", separation = "low")
y <- simulate$y
conds <- simulate$conditions
w <- rowSums(y)           ## Estimate of w
r <- table(conds)        ## Number of replicates per condition
d <- length(unique(conds)) ## Number of conditions
s <- colSums(y) / sum(y)  ## TC estimate of lib size
s.dot <- rep(NA, d)      ## Summing lib size within conditions
for(j in 1:d) s.dot[j] <- sum(s[which(conds == unique(conds)[j])]);

## Initial guess for pi and lambda
g.true <- 4
pi.guess <- simulate$pi
## Recalibrate so that (s.dot * lambda.guess) = 1
lambda.sim <- simulate$lambda
lambda.guess <- matrix(NA, nrow = d, ncol = g.true)
for(k in 1:g.true) {
  tmp <- lambda.sim[,k]/sum(lambda.sim[,k])
  lambda.guess[,k] <- tmp/s.dot
}

## Run the PMM-II model for g = 4
## with EM algorithm and "TC" library size parameter
run <- PoisMixClus(y, g = 4, norm = "TC", conds = conds)
pi.est <- run$pi
lambda.est <- run$lambda

## Mean values for each of the parameter sets
```

```

mean.guess <- PoisMixMean(y, 4, conds, s, lambda.guess)
mean.est <- PoisMixMean(y, 4, conds, s, lambda.est)

## Difference in log likelihoods
LL.diff <- logLikePoisMixDiff(y, mean.guess, pi.guess, mean.est, pi.est)
LL.diff          ## -12841.11

```

---

plot.HTSCluster	<i>Visualize results from clustering using a Poisson mixture model</i>
-----------------	--

---

## Description

A function to visualize the clustering results obtained from a Poisson mixture model.

## Usage

```

## S3 method for class 'HTSCluster'
plot(x, file.name = FALSE,
     graphs = c("map", "map.bycluster", "lambda"), data=NA, ...)
## S3 method for class 'HTSClusterWrapper'
plot(x, file.name = FALSE,
     graphs = c("capushe", "ICL", "BIC"), capushe.validation=NA, ...)

```

## Arguments

x	An object of class "HTSCluster" or "HTSClusterWrapper"
file.name	Optional file name if plots are to be saved in a PDF file.
graphs	Type of graph to be included in plots. May be equal to "map", "map.bycluster", "weighted.histograms", and/or "lambda" for objects of class "HTSCluster" and c("ICL", "BIC") for objects of class "HTSClusterWrapper"
capushe.validation	Optional number of clusters to use for capushe validation (should be less than the maximum number of clusters specified in the "HTSClusterWrapper" object).
data	( $n \times q$ ) matrix of observed counts for $n$ observations and $q$ variables (only required for the plotting of weighted histograms)
...	Additional arguments (mainly useful for plotting)

## Details

For objects of class "HTSCluster", the plotting function provides the possibility for the following visualizations:

- 1) A histogram of maximum conditional probabilities across all clusters.
- 2) Per-cluster boxplots of maximum conditional probabilities.

3) Weighted histograms of observation profiles (with weights equal to the corresponding conditional probability for each observation in each cluster), plotted independently for each variable. Fitted densities after fitting the Poisson mixture model are overlaid in red.

4) A global view of  $\lambda$  and  $\pi$  values for the selected model. When the number of conditions  $\leq 2$ , bar heights represent the value of  $\lambda_k$  for each cluster, and bar width corresponds to the value of  $\pi_k$ .

For objects of class "HTSclusterWrapper", the plotting function provides the possibility for one or all of the following visualizations:

- 1) ICL plot for all fitted models.
- 2) BIC plot for all fitted models.
- 5) Capushe diagnostic plots.

### Author(s)

Andrea Rau

### References

Rau, A., Maugis-Rabusseau, C., Martin-Magniette, M.-L., Celeux G. (2015). Co-expression analysis of high-throughput transcriptome sequencing data with Poisson mixture models. *Bioinformatics*, 31(9):1420-1427.

Andrea Rau, Gilles Celeux, Marie-Laure Martin-Magniette, and Cathy Maugis-Rabusseau (2011). Clustering high-throughput sequencing data with Poisson mixture models. *Technical report RR-7786*, Inria Saclay – Ile-de-France.

### See Also

[PoisMixClus](#), [PoisMixClusWrapper](#)

### Examples

```
set.seed(12345)

## Simulate data as shown in Rau et al. (2011)
## Library size setting "A", high cluster separation
## n = 2000 observations
simulate <- PoisMixSim(n = 200, libsize = "A", separation = "high")
y <- simulate$y
conds <- simulate$conditions

## Run the PMM-II model for g = 3
## "TC" library size estimate, EM algorithm
run <- PoisMixClus(y, g = 3,
  norm = "TC", conds = conds, init.type = "small-em")

## Visualization of results (not run):
## plot(run)
```

**Description**

These functions implement the EM and CEM algorithms for parameter estimation in a Poisson mixture model for clustering high throughput sequencing observations (e.g., genes) for a single number of clusters (PoisMixClus) or a sequence of cluster numbers (PoisMixClusWrapper). Parameters are initialized using a Small-EM strategy as described in Rau et al. (2011) or the splitting small-EM strategy described in Papastamoulis et al. (2014), and model selection is performed using the ICL criteria. Note that these functions implement the PMM-I and PMM-II models described in Rau et al. (2011).

**Usage**

```
PoisMixClus(y, g, conds, norm = "TMM",
  init.type = "small-em", init.runs = 1, init.iter = 10,
  alg.type = "EM", cutoff = 10e-6, iter = 1000, fixed.lambda = NA,
  equal.proportions = FALSE, prev.labels = NA,
  prev.probaPost = NA, verbose = FALSE, interpretation = "sum",
  EM.verbose = FALSE, wrapper = FALSE, subset.index = NA)
```

```
PoisMixClusWrapper(y, gmin = 1, gmax, conds,
  norm = "TMM", gmin.init.type = "small-em",
  init.runs = 1, init.iter = 10, split.init = TRUE, alg.type = "EM",
  cutoff = 10e-6, iter = 1000, fixed.lambda = NA,
  equal.proportions = FALSE, verbose = FALSE, interpretation = "sum",
  EM.verbose = FALSE, subset.index = NA)
```

**Arguments**

<code>y</code>	$(n \times q)$ matrix of observed counts for $n$ observations and $q$ variables
<code>g</code>	Number of clusters (a single value). If <code>fixed.lambda</code> contains a list of lambda values to be fixed, <code>g</code> corresponds to the number of clusters in addition to those fixed.
<code>gmin</code>	The minimum number of clusters in a sequence to be tested. In cases where clusters are included with a fixed value of lambda, <code>gmin</code> corresponds to the minimum number of clusters in addition to those that are fixed.
<code>gmax</code>	The maximum number of clusters in a sequence to be tested. In cases where clusters are included with a fixed value of lambda, <code>gmax</code> corresponds to the maximum number of clusters in addition to those that are fixed.
<code>conds</code>	Vector of length $q$ defining the condition (treatment group) for each variable (column) in <code>y</code>
<code>norm</code>	The type of estimator to be used to normalize for differences in library size: ("TC" for total count, "UQ" for upper quantile, "Med" for median, "DESeq" for

	the normalization method in the DESeq package, and “TMM” for the TMM normalization method (Robinson and Oshlack, 2010). Can also be a vector (of length $q$ ) containing pre-estimated library size estimates for each sample. Note that if the user provides pre-calculated normalization factors, the package will make use of <code>norm/sum(norm)</code> as normalization factors.
<code>init.type</code>	Type of initialization strategy to be used (“small-em” for the Small-EM strategy described in Rau et al. (2011), and “kmeans” for a simple $K$ -means initialization)
<code>gmin.init.type</code>	Type of initialization strategy to be used for the minimum number of clusters in a sequence ( <code>gmin</code> ): (“small-em” for the Small-EM strategy described in Rau et al. (2011), and “kmeans” for a simple $K$ -means initialization)
<code>init.runs</code>	Number of runs to be used for the Small-EM strategy described in Rau et al. (2011), with a default value of 1
<code>init.iter</code>	Number of iterations to be used within each run for the Small-EM strategy, with a default value of 10
<code>split.init</code>	If TRUE, the splitting initialization strategy of Papastamoulis et al. (2014) will be used for cluster sizes ( <code>gmin+1</code> , ..., <code>gmax</code> ). If FALSE, the initialization strategy specified in <code>gmin.init.type</code> is used for all cluster sizes in the sequence.
<code>alg.type</code>	Algorithm to be used for parameter estimation (“EM” or “CEM”)
<code>cutoff</code>	Cutoff to declare algorithm convergence (in terms of differences in log likelihoods from one iteration to the next)
<code>iter</code>	Maximum number of iterations to be run for the chosen algorithm
<code>fixed.lambda</code>	If one (or more) clusters with fixed values of <code>lambda</code> is desired, a list containing vectors of length $d$ (the number of conditions). specifying the fixed values of <code>lambda</code> for each fixed cluster.
<code>equal.proportions</code>	If TRUE, the cluster proportions are set to be equal for all clusters. Default is FALSE (unequal cluster proportions).
<code>prev.labels</code>	A vector of length $n$ of cluster labels obtained from the previous run ( $g-1$ clusters) to be used with the splitting small-EM strategy described in described in Papastamoulis et al. (2014). For other initialization strategies, this parameter takes the value NA
<code>prev.probaPost</code>	An $n \times (g-1)$ matrix of the conditional probabilities of each observation belonging to each of the $g-1$ clusters from the previous run, to be used with the splitting small-EM strategy of described in Papastamoulis et al. (2012). For other initialization strategies, this parameter takes the value NA
<code>verbose</code>	If TRUE, include verbose output
<code>interpretation</code>	If “sum”, cluster behavior is interpreted with respect to overall gene expression level (sums per gene), otherwise for “mean”, cluster behavior is interpreted with respect to mean gene expression (means per gene).
<code>EM.verbose</code>	If TRUE, more informative output is printed about the EM algorithm, including the number of iterations run and the difference between log-likelihoods at the last and penultimate iterations.
<code>subset.index</code>	Optional vector providing the indices of a subset of genes that should be used for the co-expression analysis (i.e., row indices of the data matrix $y$ ).

wrapper TRUE if the PoisMixClus function is run from within the PoisMixClusWrapper main function, and FALSE otherwise. This mainly helps to avoid recalculating parameters several times that are used throughout the algorithm (e.g., library sizes, etc.)

## Details

Output of PoisMixClus is an S3 object of class HTScluster, and output of PoisMixClusWrapper is an S3 object of class HTSclusterWrapper.

In a Poisson mixture model, the data  $\mathbf{y}$  are assumed to come from  $g$  distinct subpopulations (clusters), each of which is modeled separately; the overall population is thus a mixture of these subpopulations. In the case of a Poisson mixture model with  $g$  components, the model may be written as

$$f(\mathbf{y}; g, \Psi_g) = \prod_{i=1}^n \sum_{k=1}^g \pi_k \prod_{j=1}^d \prod_{l=1}^{r_j} P(y_{ijl}; \theta_k)$$

for  $i = 1, \dots, n$  observations in  $l = 1, \dots, r_j$  replicates of  $j = 1, \dots, d$  conditions (treatment groups), where  $P(\cdot)$  is the standard Poisson density,  $\Psi_g = (\pi_1, \dots, \pi_{g-1}, \theta')$ ,  $\theta'$  contains all of the parameters in  $\theta_1, \dots, \theta_g$  assumed to be distinct, and  $\pi = (\pi_1, \dots, \pi_g)'$  are the mixing proportions such that  $\pi_k$  is in  $(0,1)$  for all  $k$  and  $\sum_k \pi_k = 1$ .

We consider the following parameterization for the mean  $\theta_k = (\mu_{ijlk})$ . We consider

$$\mu_{ijlk} = w_i s_{jl} \lambda_{jk}$$

where  $w_i$  corresponds to the expression level of observation  $i$ ,  $\lambda_k = (\lambda_{1k}, \dots, \lambda_{dk})$  corresponds to the clustering parameters that define the profiles of the genes in cluster  $k$  across all variables, and  $s_{jl}$  is the normalized library size (a fixed constant) for replicate  $l$  of condition  $j$ .

There are two approaches to estimating the parameters of a finite mixture model and obtaining a clustering of the data: the estimation approach (via the EM algorithm) and the clustering approach (via the CEM algorithm). Parameter initialization is done using a Small-EM strategy as described in Rau et al. (2011) via the `emInit` function. Model selection may be performed using the BIC or ICL criteria, or the slope heuristics.

## Value

lambda	$(d \times g)$ matrix containing the estimate of $\hat{\lambda}$
pi	Vector of length $g$ containing the estimate of $\hat{\pi}$
labels	Vector of length $n$ containing the cluster assignments of the $n$ observations
probaPost	Matrix containing the conditional probabilities of belonging to each cluster for all observations
log.like	Value of log likelihood
BIC	Value of BIC criterion
ICL	Value of ICL criterion
alg.type	Estimation algorithm used; matches the argument <code>alg.type</code> above)

<code>norm</code>	Library size normalization factors used
<code>conds</code>	Conditions specified by user
<code>iterations</code>	Number of iterations run
<code>logLikeDiff</code>	Difference in log-likelihood between the last and penultimate iterations of the algorithm
<code>subset.index</code>	If provided by the user, the indices of subset of genes used for co-expression analyses
<code>loglike.all</code>	Log likelihoods calculated for each of the fitted models for cluster sizes <code>gmin</code> , ..., <code>gmax</code>
<code>capushe</code>	Results of capushe model selection, an object of class "Capushe"
<code>ICL.all</code>	ICL values calculated for each of the fitted models for cluster sizes <code>gmin</code> , ..., <code>gmax</code>
<code>ICL.results</code>	Object of class HTScluster giving the results from the model chosen via the ICL criterion
<code>BIC.results</code>	Object of class HTScluster giving the results from the model chosen via the BIC
<code>DDSE.results</code>	Object of class HTScluster giving the results from the model chosen via the DDSE slope heuristics criterion
<code>Djump.results</code>	Object of class HTScluster giving the results from the model chosen via the Djump slope heuristics criterion
<code>all.results</code>	List of objects of class HTScluster giving the results for all models for cluster sizes <code>gmin</code> , ..., <code>gmax</code>
<code>model.selection</code>	Type of criteria used for model selection, equal to NA for direct calls to <code>PoisMixClus</code> or "DDSE", "Djump", "BIC", or "ICL" for the respective selected models for calls to <code>PoisMixClusWrapper</code>

### Note

Note that the `fixed.lambda` argument is primarily intended to be used in the case when a single cluster is fixed to have equal clustering parameters lambda across all conditions (i.e.,  $\lambda_{j1} = \lambda_1 = 1$ ); this is particularly useful when identifying genes with non-differential expression across all conditions (see the HTSDiff R package for more details). Alternatively, this argument could be used to specify a cluster for which genes are only expressed in a single condition (e.g.,  $\lambda_{11} = 1$  and  $\lambda_{j1} = 0$  for all  $j > 1$ ). Other possibilities could be considered, but note that the fixed values of lambda must satisfy the constraint  $\sum_j \lambda_{jk} s_{j.} = 1$  for all  $k$  imposed in the model; if this is not the case, a warning message will be printed.

### Author(s)

Andrea Rau



## References

- Anders, S. and Huber, W. (2010) Differential expression analysis for sequence count data. *Genome Biology*, **11**(R106), 1-28.
- Papastamoulis, P., Martin-Magniette, M.-L., and Maugis-Rabuseau, C. (2014). On the estimation of mixtures of Poisson regression models with large number of components. *Computational Statistics and Data Analysis*: 3rd special Issue on Advances in Mixture Models, DOI: 10.1016/j.csda.2014.07.005.
- Rau, A., Maugis-Rabuseau, C., Martin-Magniette, M.-L., Celeux G. (2015). Co-expression analysis of high-throughput transcriptome sequencing data with Poisson mixture models. *Bioinformatics*, **31**(9):1420-1427.
- Rau, A., Celeux, G., Martin-Magniette, M.-L., Maugis-Rabuseau, C (2011). Clustering high-throughput sequencing data with Poisson mixture models. Inria Research Report 7786. Available at <https://inria.hal.science/inria-00638082>.

## See Also

[probaPost](#) for the calculation of the conditional probability of belonging to a cluster; [PoisMixMean](#) for the calculation of the per-cluster conditional mean of each observation; [logLikePoisMixDiff](#) for the calculation of the log likelihood of a Poisson mixture model; [emInit](#) and [kmeanInit](#) for the Small-EM parameter initialization strategy

## Examples

```
set.seed(12345)

## Simulate data as shown in Rau et al. (2011)
## Library size setting "A", high cluster separation
## n = 200 observations

simulate <- PoisMixSim(n = 200, libsize = "A", separation = "high")
y <- simulate$y
conds <- simulate$conditions

## Run the PMM model for g = 3
## "TC" library size estimate, EM algorithm

run <- PoisMixClus(y, g = 3, conds = conds, norm = "TC")

## Estimates of pi and lambda for the selected model

pi.est <- run$pi
lambda.est <- run$lambda

## Not run: PMM for 4 total clusters, with one fixed class
## "TC" library size estimate, EM algorithm
##
## run <- PoisMixClus(y, g = 3, norm = "TC", conds = conds,
##   fixed.lambda = list(c(1,1,1)))
##
##
```

```

## Not run: PMM model for 4 clusters, with equal proportions
## "TC" library size estimate, EM algorithm
##
## run <- PoisMixClus(y, g = 4, norm = "TC", conds = conds,
##   equal.proportions = TRUE)
##
##
## Not run: PMM model for g = 1, ..., 10 clusters, Split Small-EM init
##
## run1.10 <- PoisMixClusWrapper(y, gmin = 1, gmax = 10, conds = conds,
## norm = "TC")
##
##
## Not run: PMM model for g = 1, ..., 10 clusters, Small-EM init
##
## run1.10bis <- <- PoisMixClusWrapper(y, gmin = 1, gmax = 10, conds = conds,
## norm = "TC", split.init = FALSE)
##
##
## Not run: previous model equivalent to the following
##
## for(K in 1:10) {
## run <- PoisMixClus(y, g = K, conds = conds, norm = "TC")
## }

```

---

PoisMixMean

*Calculate the conditional per-cluster mean of each observation*


---

### Description

This function is used to calculate the conditional per-cluster mean expression for all observations. This value corresponds to  $\mu = (\mu_{ijkl}) = (\hat{w}_i \hat{\lambda}_{jk})$  for the PMM-I model and  $\mu = (\mu_{ijkl}) = (\hat{w}_i s_{jl} \hat{\lambda}_{jk})$  for the PMM-II model.

### Usage

```
PoisMixMean(y, g, conds, s, lambda)
```

### Arguments

y	( $n \times q$ ) matrix of observed counts for $n$ observations and $q$ variables
g	Number of clusters
conds	Vector of length $q$ defining the condition (treatment group) for each variable (column) in y
s	Estimate of normalized per-variable library size
lambda	( $d \times g$ ) matrix containing the current estimate of lambda, where $d$ is the number of conditions (treatment groups) and $g$ is the number of clusters

**Value**

A list of length  $g$  containing the  $(n \times q)$  matrices of mean expression for all observations, conditioned on each of the  $g$  clusters

**Author(s)**

Andrea Rau

**References**

Rau, A., Maugis-Rabusseau, C., Martin-Magniette, M.-L., Celeux G. (2015). Co-expression analysis of high-throughput transcriptome sequencing data with Poisson mixture models. *Bioinformatics*, 31(9):1420-1427.

Rau, A., Celeux, G., Martin-Magniette, M.-L., Maugis-Rabusseau, C. (2011). Clustering high-throughput sequencing data with Poisson mixture models. Inria Research Report 7786. Available at <https://inria.hal.science/inria-00638082>.

**See Also**

[PoisMixClus](#) for Poisson mixture model estimation and model selection

**Examples**

```
set.seed(12345)

## Simulate data as shown in Rau et al. (2011)
## Library size setting "A", high cluster separation
## n = 200 observations

simulate <- PoisMixSim(n = 200, libsize = "A", separation = "high")
y <- simulate$y
conds <- simulate$conditions
s <- colSums(y) / sum(y) ## TC estimate of lib size

## Run the PMM-II model for g = 3
## "TC" library size estimate, EM algorithm

run <- PoisMixClus(y, g = 3, norm = "TC", conds = conds)
pi.est <- run$pi
lambda.est <- run$lambda

## Calculate the per-cluster mean for each observation
means <- PoisMixMean(y, g = 3, conds, s, lambda.est)
```

PoisMixSim

*Simulate data from a Poisson mixture model***Description**

This function simulates data from a Poisson mixture model, as described by Rau et al. (2011). Data are simulated with varying expression level ( $w_i$ ) for 4 clusters. Clusters may be simulated with “high” or “low” separation, and three different options are available for the library size setting: “equal”, “A”, and “B”, as described by Rau et al. (2011).

**Usage**

```
PoisMixSim(n = 2000, libsize, separation)
```

**Arguments**

n	Number of observations
libsize	The type of library size difference to be simulated (“equal”, “A”, or “B”, as described by Rau et al. (2011))
separation	Cluster separation (“high” or “low”, as described by Rau et al. (2011))

**Value**

y	( $n \times q$ ) matrix of simulated counts for $n$ observations and $q$ variables
labels	Vector of length $n$ defining the true cluster labels of the simulated data
pi	Vector of length 4 (the number of clusters) containing the true value of $\pi$
lambda	( $d \times 4$ ) matrix of $\lambda$ values for $d$ conditions (3 in the case of libsize = “equal” or “A”, and 2 otherwise) in 4 clusters (see note below)
w	Row sums of y (estimate of $\hat{w}$ )
conditions	Vector of length $q$ defining the condition (treatment group) for each variable (column) in y

**Note**

If one or more observations are simulated such that all variables have a value of 0, those rows are removed from the data matrix; as such, in some cases the simulated data y may have less than n rows.

The PMM-I model includes the parameter constraint  $\sum_k \lambda_{jk} r_j = 1$ , where  $r_j$  is the number of replicates in condition (treatment group)  $j$ . Similarly, the parameter constraint in the PMM-II model is  $\sum_j \sum_l \lambda_{jk} s_{jl} = 1$ , where  $s_{jl}$  is the library size for replicate  $l$  of condition  $j$ . The value of lambda corresponds to that used to generate the simulated data, where the library sizes were set as described in Table 2 of Rau et al. (2011). However, due to variability in the simulation process, the actually library sizes of the data y are not exactly equal to these values; this means that the value of lambda may not be directly compared to an estimated value of  $\hat{\lambda}$  as obtained from the [PoisMixClus](#) function.

**Author(s)**

Andrea Rau

**References**

Rau, A., Celeux, G., Martin-Magniette, M.-L., Maugis-Rabusseau, C. (2011). Clustering high-throughput sequencing data with Poisson mixture models. Inria Research Report 7786. Available at <https://inria.hal.science/inria-00638082>.

**Examples**

```
set.seed(12345)

## Simulate data as shown in Rau et al. (2011)
## Library size setting "A", high cluster separation
## n = 200 observations

simulate <- PoisMixSim(n = 200, libsize = "A", separation = "high")
y <- simulate$y
conds <- simulate$conditions
```

---

 probaPost

---

*Calculate the conditional probability of belonging to each cluster in a Poisson mixture model*


---

**Description**

This function computes the conditional probabilities  $t_{ik}$  that an observation  $i$  arises from the  $k^{\text{th}}$  component for the current value of the mixture parameters.

**Usage**

```
probaPost(y, g, conds, pi, s, lambda)
```

**Arguments**

y	( $n \times q$ ) matrix of observed counts for $n$ observations and $q$ variables
g	Number of clusters
conds	Vector of length $q$ defining the condition (treatment group) for each variable (column) in y
pi	Vector of length $g$ containing the current estimate of $\hat{\pi}$
s	Vector of length $q$ containing the estimates for the normalized library size parameters for each of the $q$ variables in y
lambda	( $d \times g$ ) matrix containing the current estimate $\lambda$ , where $d$ is the number of conditions (treatment groups)

**Value**

t  $(n \times g)$  matrix made up of the conditional probability of each observation belonging to each of the  $g$  clusters

**Note**

If all values of  $t_{ik}$  are 0 (or nearly zero), the observation is assigned with probability one to belong to the cluster with the closest mean (in terms of the Euclidean distance from the observation). To avoid calculation difficulties, extreme values of  $t_{ik}$  are smoothed, such that those smaller than  $1e-10$  or larger than  $1-1e-10$  are set equal to  $1e-10$  and  $1-1e-10$ , respectively.

**Author(s)**

Andrea Rau

**References**

Rau, A., Maugis-Rabusseau, C., Martin-Magniette, M.-L., Celeux G. (2015). Co-expression analysis of high-throughput transcriptome sequencing data with Poisson mixture models. *Bioinformatics*, 31(9):1420-1427.

Rau, A., Celeux, G., Martin-Magniette, M.-L., Maugis-Rabusseau, C. (2011). Clustering high-throughput sequencing data with Poisson mixture models. Inria Research Report 7786. Available at <https://inria.hal.science/inria-00638082>.

**See Also**

[PoisMixClus](#) for Poisson mixture model estimation and model selection; [PoisMixMean](#) to calculate the conditional per-cluster mean of each observation

**Examples**

```
set.seed(12345)

## Simulate data as shown in Rau et al. (2011)
## Library size setting "A", high cluster separation
## n = 200 observations

simulate <- PoisMixSim(n = 200, libsize = "A", separation = "high")
y <- simulate$y
conds <- simulate$conditions
s <- colSums(y) / sum(y)    ## TC estimate of lib size

## Run the PMM-II model for g = 3
## "TC" library size estimate, EM algorithm

run <- PoisMixClus(y, g = 3, norm = "TC",
  conds = conds)
pi.est <- run$pi
lambda.est <- run$lambda
```

```
## Calculate the conditional probability of belonging to each cluster
proba <- probaPost(y, g = 3, conds = conds, pi = pi.est, s = s,
lambda = lambda.est)

## head(round(proba,2))
```

---

summary.HTSCluster      *Summarize results from clustering using a Poisson mixture model*

---

### Description

A function to summarize the clustering results obtained from a Poisson mixture model.

### Usage

```
## S3 method for class 'HTSCluster'
summary(object, ...)
## S3 method for class 'HTSClusterWrapper'
summary(object, ...)
```

### Arguments

object	An object of class "HTSCluster" or "HTSClusterWrapper"
...	Additional arguments

### Details

The summary function for an object of class "HTSCluster" provides the following summary of results:

- 1) Number of clusters and model selection criterion used, if applicable.
- 2) Number of observations across all clusters with a maximum conditional probability greater than 90 model.
- 3) Number of observations per cluster with a maximum conditional probability greater than 90 selected model.
- 4)  $\lambda$  values for the selected model.
- 5)  $\pi$  values for the selected model.

The summary function for an object of class "HTSClusterWrapper" provides the number of clusters selected for the BIC, ICL, DDSE, and Djump model selection approaches.

### Author(s)

Andrea Rau

## References

Rau, A., Maugis-Rabusseau, C., Martin-Magniette, M.-L., Celeux G. (2015). Co-expression analysis of high-throughput transcriptome sequencing data with Poisson mixture models. *Bioinformatics*, 31(9):1420-1427.

Rau, A., Celeux, G., Martin-Magniette, M.-L., Maugis-Rabusseau, C. (2011). Clustering high-throughput sequencing data with Poisson mixture models. Inria Research Report 7786. Available at <https://inria.hal.science/inria-00638082>.

## See Also

[PoisMixClus](#), [PoisMixClusWrapper](#)

## Examples

```
set.seed(12345)

## Simulate data as shown in Rau et al. (2011)
## Library size setting "A", high cluster separation
## n = 2000 observations
simulate <- PoisMixSim(n = 200, libsize = "A", separation = "high")
y <- simulate$y
conds <- simulate$conditions

## Run the PMM-II model for g = 3
## "TC" library size estimate, EM algorithm
run <- PoisMixClus(y, g = 3,
  norm = "TC", conds = conds, init.type = "small-em")

## Summary of results:
summary(run)
```



# Index

- \* **cluster**
  - HTSCluster-package, 2
  - PoisMixClus, 13
- \* **datagen**
  - PoisMixSim, 20
- \* **documentation**
  - HTSClusterUsersGuide, 4
- \* **methods**
  - highDimensionARI, 3
  - logLikePoisMix, 8
  - plot.HTSCluster, 11
  - PoisMixMean, 18
  - probaPost, 21
  - summary.HTSCluster, 23
- \* **models**
  - HTSCluster-package, 2
  - Init, 5
  - plot.HTSCluster, 11
  - PoisMixClus, 13

emInit, 15, 17

emInit (Init), 5

highDimensionARI, 3

HTSCluster (HTSCluster-package), 2

HTSCluster-package, 2

HTSClusterUsersGuide, 4

Init, 5

kmeanInit, 17

kmeanInit (Init), 5

logLikePoisMix, 8

logLikePoisMixDiff, 17

logLikePoisMixDiff (logLikePoisMix), 8

mylogLikePoisMixObs (logLikePoisMix), 8

plot.HTSCluster, 11

plot.HTSClusterWrapper

- (plot.HTSCluster), 11

PoisMixClus, 8–10, 12, 13, 19, 20, 22, 24

PoisMixClusWrapper, 8, 12, 24

PoisMixClusWrapper (PoisMixClus), 13

PoisMixMean, 9, 10, 17, 18, 22

PoisMixSim, 20

probaPost, 17, 21

probaPostInit (Init), 5

splitEMInit (Init), 5

summary.HTSCluster, 23

summary.HTSClusterWrapper

- (summary.HTSCluster), 23

Sweave, 4

system, 5