

# Package ‘essHist’

May 10, 2019

**Type** Package

**Title** The Essential Histogram

**Version** 1.2.2

**Date** 2019-05-10

**Author** Housen Li [aut, cre],  
Hannes Sieling [aut]

**Maintainer** Housen Li <housen.li@outlook.com>

**Description** Provide an optimal histogram, in the sense of probability density estimation and features detection, by means of multiscale variational inference. In other words, the resulting histogram serves as an optimal density estimator, and meanwhile recovers the features, such as increases or modes, with both false positive and false negative controls. Moreover, it provides a parsimonious representation in terms of the number of blocks, which simplifies data interpretation. The only assumption for the method is that data points are independent and identically distributed, so it applies to fairly general situations, including continuous distributions, discrete distributions, and mixtures of both. For details see Li, Munk, Sieling and Walther (2016) <arXiv:1612.07216>.

**Depends** R (>= 2.15.3)

**License** GPL-3

**LazyData** TRUE

**Imports** Rcpp (>= 0.12.5), graphics, stats, grDevices, utils

**LinkingTo** Rcpp

**Suggests** testthat

**NeedsCompilation** yes

**Repository** CRAN

**Date/Publication** 2019-05-10 09:30:03 UTC

## R topics documented:

essHist-package . . . . .	2
checkHistogram . . . . .	3
Essential Histogram . . . . .	6

Generate Intervals . . . . .	8
Mixed normals . . . . .	9
Multiscale Quantiles . . . . .	10

<b>Index</b>	<b>13</b>
--------------	-----------

---

essHist-package	<i>The Essential Histogram</i>
-----------------	--------------------------------

---

## Description

Provide an optimal histogram, in the sense of probability density estimation and features detection, by means of multiscale variational inference. In other words, the resulting histogram serves as an optimal density estimator, and meanwhile recovers the features, such as increases or modes, with both false positive and false negative controls. Moreover, it provides a parsimonious representation in terms of the number of blocks, which simplifies data interpretation. The only assumption for the method is that data points are independent and identically distributed, so it applies to fairly general situations, including continuous distributions, discrete distributions, and mixtures of both. For details see Li, Munk, Sieling and Walther (2016) <arXiv:1612.07216>.

## Details

Package: essHist  
 Type: Package  
 Version: 1.2.2  
 Date: 2019-05-10  
 License: GPL-3

### Index:

essHistogram	Compute the essential histogram
checkHistogram	Check any estimator by the multiscale confidence set
genIntv	Generate the system of intervals
msQuantile	Simulate the quantile of multiscale statistics
dmixnorm	Compute density function of Gaussian mixtures
pmixnorm	Compute distribution function of Gaussian mixtures
rmixnorm	Generate random number of Gaussian mixtures
paramExample	Output detailed parameters for some famous examples

## Author(s)

Housen Li [aut, cre], Hannes Sieling [aut]

Maintainer: Housen Li <housen.li@outlook.com>

## References

Li, H., Munk, A., Sieling, H., and Walther, G. (2016). The essential histogram. arXiv:1612.07216

## Examples

```
# Simulate data
set.seed(123)
type = 'skewed_unimodal'
n = 500
y = rmixnorm(n, type = type)

# Compute the essential histogram
eh = essHistogram(y, plot = FALSE)

# Plot results
#   compute oracle density
x = sort(y)
od = dmixnorm(x, type = type)
#   compare with orcle density
plot(x, od, type = "l", xlab = NA, ylab = NA, col = "red", main = type)
lines(eh)
legend("topleft", c("Oracle density", "Essential histogram"),
      lty = c(1,1), col = c("red", "black"))

##### Evaluate other method
set.seed(123)
# Data: mixture of Gaussians "harp"
n = 500
y = rmixnorm(n, type = 'harp')

# Oracle density
x = sort(y)
ho = dmixnorm(x, type = 'harp')

# R default histogram
h = hist(y, plot = FALSE)

# Check R default histogram to local multiscale constraints
b = checkHistogram(h, y, ylim=c(-0.1,0.16))
lines(x, ho, col = "red")
rug(x, col = 'blue')
legend("topright", c("R-Histogram", "Truth"), col = c("black", "red"), lty = c(1,1))
```

---

checkHistogram

*Check any histogram estimator by means of the multiscale confidence set*

---

## Description

Provide the locations, i.e., intervals, where features are potentially missing (a.k.a. false negatives), and the break-points that are potentially redundant (a.k.a. false positives), by means of the multiscale confidence set.

## Usage

```
checkHistogram(h, x, alpha = 0.1, q = NULL, intv = NULL,
              mode = ifelse(anyDuplicated(x), "Gen", "Con"),
              plot = TRUE, xlim = NULL, ylim = NULL,
              xlab = "", ylab = "", yaxt = "n", ...)
```

## Arguments

h	a numeric vector specifying values of a histogram at sample points; or a histogram class object (i.e. the return value of <a href="#">hist</a> ).
x	a numeric vector containing the data.
alpha	significance level, default as 0.1, see also <a href="#">essHistogram</a> .
q	threshold of the multiscale constraint; by default, q is chosen as the (1-alpha)-quantile of the null distribution of the multiscale statistic via Monte Carlo simulation, see also <a href="#">msQuantile</a> .
intv	a data frame provides the system of intervals on which the multiscale statistic is defined. The data frame contains the following two columns left left index of an interval right right index of an interval By default, it is set to the sparse interval system proposed by Rivera and Walther (2013), see also Li et al. (2016).
mode	"Con" for continuous distribution functions "Gen" for general (possibly with discontinuous) distribution functions By default, "Con" is chosen if there is no tied observations; otherwise, "Gen" is chosen; see Li et al. (2016) for further details.
plot	logical. If TRUE, the input estimator is plotted, together with evaluation information. More precisely, at the very bottom, intervals where local constraints are violated are plotted. In the middle short vertical lines that indicate possibly removable change-points are drawn above a light blue horizontal line. Right below the light blue line, it plots a horizontal gray scale strap, the darkness of which reflects the number of violation intervals covering a given location, as a summary of violation information.
xlim, ylim	numeric vectors of length 2 (default xlim = range(y), ylim = NULL): see <a href="#">plot</a> .
xlab	a title for the x axis (default empty string): see <a href="#">title</a> and <a href="#">plot</a> .
ylab	a title for the y axis (default empty string): see <a href="#">title</a> and <a href="#">plot</a> .
yaxt	A character which specifies the y axis type (default "n"): see <a href="#">par</a> .
...	further arguments and <a href="#">graphical parameters</a> passed to <a href="#">plot</a> (if plot = TRUE).

**Details**

This function presents a visualization: the upper part plots the given histogram; in the middle part short vertical lines mark all removable break-points; in the lower part intervals of violation are shown, and a graybar below the middle horizontal line (blue) summarizes such violations with the darkness scaling with the number of violation intervals covering a location. See Examples below and Li et al. (2016) for further details.

**Value**

A list consists of one data frame, and one numeric vector:

violatedIntervals

A data frame provides the intervals where the corresponding local side constraint is violated; an empty data frame if there is no violation. It contains the following four columns

leftIndex left index of an interval

rightIndex right index of an interval

leftEnd left end point of an interval

rightEnd right end point of an interval

An empty data.frame is returned if there is no violation.

removableBreakpoints

A numeric vector contains all removable breakpoints, with zero length if there is no removable breakpoint.

**Note**

The argument `intv` is internally adjusted ensure it contains no empty intervals in case of tied observations. Only the intervals on which the input histogram is constant will be checked! All the printing messages can be disabled by calling `suppressMessages`.

**References**

Li, H., Munk, A., Sieling, H., and Walther, G. (2016). The essential histogram. arXiv:1612.07216.

**See Also**

[essHistogram](#), [genIntv](#), [msQuantile](#)

**Examples**

```
set.seed(123)
# Data: mixture of Gaussians "harp"
n = 500
y = rmixnorm(n, type = 'harp')

# Oracle density
x = sort(y)
ho = dmixnorm(x, type = 'harp')
```

```
# R default histogram
h = hist(y, plot = FALSE)

# Check R default histogram to local multiscale constraints
b = checkHistogram(h, y, ylim=c(-0.1,0.16))
lines(x, ho, col = "red")
rug(x, col = 'blue')
legend("topright", c("R-Histogram", "Truth"), col = c("black", "red"), lty = c(1,1))
```

---

Essential Histogram     *The Essential Histogram*

---

## Description

Compute the essential histogram via (pruned) dynamic programming.

## Usage

```
essHistogram(x, alpha = 0.5, q = NULL, intv = NULL, plot = TRUE,
             mode = ifelse(anyDuplicated(x), "Gen", "Con"),
             xname = deparse(substitute(x)), ...)
```

## Arguments

<code>x</code>	a numeric vector containing the data.
<code>alpha</code>	significance level; default as 0.5. One should set <code>alpha = 0.1</code> or even smaller if confidence statements have to be made, while one can set <code>alpha = 0.9</code> if the goal is to explore the data for potential features with tolerance to false positives. The default value is only a trade-off.
<code>q</code>	threshold value; by default, <code>q</code> is chosen as the $(1-\alpha)$ -quantile of the null distribution of the multiscale statistic via Monte Carlo simulation, see also <a href="#">msQuantile</a> .
<code>intv</code>	a data frame provides the system of intervals on which the multiscale statistic is defined. The data frame contains the following two columns <code>left</code> left index of an interval <code>right</code> right index of an interval By default, it is set to the sparse interval system proposed by Rivera and Walther (2013), see also Li et al. (2016).
<code>plot</code>	logical. If TRUE (default), a histogram is plotted, otherwise a list of breaks and counts is returned. In the latter case, a warning is used if (typically graphical) arguments are specified that only apply to the <code>plot = TRUE</code> case.
<code>mode</code>	"Con" for continuous distribution functions "Gen" for general (possibly with discontinuous) distribution functions By default, "Con" is chosen if there is no tied observations; otherwise, "Gen" is chosen; see Li et al. (2016) for further details.
<code>xname</code>	a character string with the actual <code>x</code> argument name.
<code>...</code>	further arguments and <a href="#">graphical parameters</a> passed to <code>plot.histogram</code> and thence to <code>title</code> and <code>axis</code> (if <code>plot = TRUE</code> ).

## Details

The essential histogram is defined as the histogram with least blocks within the multiscale constraint. The one with highest likelihood is picked if there are more than one solutions. The essential histogram involves only one parameter  $q$ , the threshold of the multiscale constraint. Such a parameter can be chosen by means of the significance level  $\alpha$ , which leads to nature statistical significance statements for the multiscale constraint. The computational complexity is often linear in terms of sample size, although the worst complexity bound is quadratic up to a log-factor in case of the sparse interval system. See Li et al. (2016) for further details.

## Value

An object of class "histogram", which is of the same class as returned by function [hist](#).

## Note

The argument `intv` is internally adjusted to ensure it contains no empty intervals, especially in case of tied observations. The first block of the returned histogram is a closed interval, and the rest blocks are left open right closed intervals. All the printing messages can be disabled by calling [suppressMessages](#).

## References

Li, H., Munk, A., Sieling, H., and Walther, G. (2016). The essential histogram. arXiv:1612.07216.  
Rivera, C., & Walther, G. (2013). Optimal detection of a jump in the intensity of a Poisson process or in a density with likelihood ratio statistics. *Scand. J. Stat.* 40, 752–769.

## See Also

[checkHistogram](#), [genIntv](#), [hist](#), [msQuantile](#)

## Examples

```
# Simulate data
set.seed(123)
type = 'skewed_unimodal'
n = 500
y = rmixnorm(n, type = type)

# Compute the essential histogram
eh = essHistogram(y, plot = FALSE)

# Plot results
#   compute oracle density
x = sort(y)
od = dmixnorm(x, type = type)
#   compare with orcle density
plot(x, od, type = "l", xlab = NA, ylab = NA, col = "red", main = type)
lines(eh)
legend("topleft", c("Oracle density", "Essential histogram"),
      lty = c(1,1), col = c("red", "black"))
```

---

Generate Intervals      *Generate the system of intervals*

---

### Description

Generate the system of intervals on which the multiscale statistic is defined, see Li et al. (2016).

### Usage

```
genIntv(n, type = c("Sparse", "Full"))
```

### Arguments

n	number of observations.
type	type of interval system. type = "Sparse" (default) is the sparse system proposed by Rivera and Walther (2013), see also Li et al. (2016). type = "Full" is the system of all possible intervals with end-index ranging from 1 to n.

### Value

A data frame provides the system of intervals, and consists two columns

left	left index of an interval
right	right index of an interval

### References

Li, H., Munk, A., Sieling, H., and Walther, G. (2016). The essential histogram. arXiv:1612.07216.  
Rivera, C., & Walther, G. (2013). Optimal detection of a jump in the intensity of a Poisson process or in a density with likelihood ratio statistics. *Scand. J. Stat.* 40, 752–769.

### See Also

[checkHistogram](#), [essHistogram](#), [msQuantile](#)

### Examples

```
n = 5  
intv = genIntv(n, "Full")  
print(intv)
```



---

Mixed normals

*The mixture of normal distributions*


---

## Description

Density, distribution function and random generation for the mixture of normals with each component specified by mean and sd, and mixture weights by prob. `paramExample` gives detailed parameters for some examples specified by type.

## Usage

```
dmixnorm(x, mean = c(0), sd = rep(1,length(mean)),
          prob = rep(1,length(mean)), type = NULL, ...)
pmixnorm(x, mean = c(0), sd = rep(1,length(mean)),
          prob = rep(1,length(mean)), type = NULL, ...)
rmixnorm(n, mean = c(0), sd = rep(1,length(mean)),
          prob = rep(1,length(mean)), type = NULL)
paramExample(type)
```

## Arguments

<code>x</code>	vector of locations.
<code>n</code>	integer; number of observations.
<code>mean</code>	vector of means for each mixture component.
<code>sd</code>	vector of standard deviations for each mixture component. Default is of unit variance for each component.
<code>prob</code>	vector of prior probability for each mixture component (i.e. mixture weights). All nonnegative values are allowed, and automatically recaled to ensure their sum equal to 1. Default is of equal probability for each component.
<code>type</code>	a (case insensitive) character string of example name; It includes examples from Marron & Wand (1992): "MW1", ..., "MW15", or equivalently "gauss", "skewed_unimodal", "strong_skewed", "kurtotic_unimodal", "outlier", "bimodal", "separated_bimodal", "skewed_bimodal", "trimodal", "claw", "double_claw", "asymmetric_claw", "asymmetric_double_claw", "smooth_comb", "discrete_comb"; It also includes "harp" example from Li et al. (2016).
<code>...</code>	further arguments passed to <code>dnorm</code> and <code>pnorm</code> .

## Details

Users either provide, optionally, mean, sd and prob; or type. In case of providing type, the values of mean, sd and prob are ignored.

The default case is standard normal, the same as `dnorm`, `pnorm` and `rnorm`.

**Value**

dmixnorm gives the density, pmixnorm gives the distribution function, and rmixnorm generates random deviates.

The length of the result is determined by n for rmixnorm, and is the length of x for dmixnorm and pmixnorm.

paramExample gives a data frame with components mean, sd and prob.

**References**

Li, H., Munk, A., Sieling, H., and Walther, G. (2016). The essential histogram. arXiv:1612.07216.

Marron, J. S., & Wand, M. P. (1992). Exact mean integrated squared error. *Ann. Statist.*, 20(2), 712–736.

**See Also**

[Normal](#) for standard normal distributions; [Distributions](#) for other standard distributions.

**Examples**

```
## Example harp
type = "harp"
#   generate random numbers
n = 500
Y = rmixnorm(n, type = type)
#   compute the density
x = seq(min(Y), max(Y), length.out = n)
f = dmixnorm(x, type = type)
#   compute the distribution
F = pmixnorm(x, type = type)
#   plots
op = par(mfrow = c(1,2))
plot(x, f, type = "l", main = "Harp Density")
rug(Y, col = 'red')
plot(x, F, type = "l", main = "Harp Distribution")
rug(Y, col = 'red')
par(op)
```

---

Multiscale Quantiles *Quantile of the multiscale statistics*

---

**Description**

Simulate quantiles of the multiscale statistics under any continuous distribution function.

**Usage**

```
msQuantile(n, alpha = c(0.5), nsim = 5e3, is.sim = (n < 1e4),
           intv = genIntv(n), mode = c("Con", "Gen"), ...)
```

**Arguments**

<code>n</code>	number of observations.
<code>alpha</code>	significance level; default as 0.5, see also <a href="#">essHistogram</a> . Like <a href="#">quantile</a> , it can also be a vector.
<code>nsim</code>	number of Monte Carlo simulations.
<code>is.sim</code>	logical. If TRUE (default if $n < 10,000$ ) the quantile is determined via Monte Carlo simulations, which might take a long time; otherwise (default if $n \geq 10,000$ ) it uses the quantile with $n = 10,000$ , which has been precomputed and stored.
<code>intv</code>	a data frame provides the system of intervals on which the multiscale statistic is defined. The data frame contains the following two columns <code>left</code> left index of an interval <code>right</code> right index of an interval By default, it is set to the sparse interval system proposed by Rivera and Walther (2013), see <a href="#">genIntv</a> and also Li et al. (2016).
<code>mode</code>	"Con" for continuous distribution functions (default) "Gen" for general (possibly with discontinuous) distribution functions See Li et al. (2016) for further details.
<code>...</code>	further arguments passed to function <a href="#">quantile</a> .

**Details**

Empirically, it turns out that the quantile of the multiscale statistic converges fast to that of the limit distribution as the number of observations  $n$  increases. Thus, for the sake of computational efficiency, the quantile with  $n = 10,000$  are used by default for that with  $n > 10,000$ , which has already been precomputed and stored. Of course, for arbitrary sample size  $n$ , one can always simulate the quantile by setting `is.sim = TRUE`, and use the precomputed value by setting `is.sim = FALSE`. For a given sample size  $n$ , simulations are once computed, and then automatically recorded in the R memory for later usage. For memory efficiency, only the last simulation is stored.

**Value**

A vector of length `length(alpha)` is returned, the same structure as returned by function [quantile](#) with option `names = FALSE`; The values are the  $(1-\alpha)$ -quantile(s) of the null distribution of the multiscale statistic via Monte Carlo simulation, corresponding to  $(1-\alpha)$ -confidence level(s). See Li et al. (2016) for further details.

**Note**

All the printing messages can be disabled by calling [suppressMessages](#).

**References**

- Li, H., Munk, A., Sieling, H., and Walther, G. (2016). The essential histogram. arXiv:1612.07216.
- Rivera, C., & Walther, G. (2013). Optimal detection of a jump in the intensity of a Poisson process or in a density with likelihood ratio statistics. *Scand. J. Stat.* 40, 752–769.

**See Also**

[checkHistogram](#), [essHistogram](#), [genIntv](#)

**Examples**

```
n = 100 # number of observations
nsim = 100 # number of simulations

alpha = c(0.1, 0.9) # significance level
q = msQuantile(n, alpha, nsim)

print(q)
```

# Index

- \*Topic **datagen**
  - Mixed normals, [9](#)
- \*Topic **distribution**
  - checkHistogram, [3](#)
  - Essential Histogram, [6](#)
  - essHist-package, [2](#)
  - Mixed normals, [9](#)
  - Multiscale Quantiles, [10](#)
- \*Topic **nonparametric**
  - Essential Histogram, [6](#)
  - essHist-package, [2](#)
  - Generate Intervals, [8](#)
  - Multiscale Quantiles, [10](#)
- \*Topic **package**
  - essHist-package, [2](#)
- axis, [6](#)
- checkHistogram, [3](#), [7](#), [8](#), [12](#)
- Distributions, [10](#)
- dmixnorm (Mixed normals), [9](#)
- dnorm, [9](#)
- Essential Histogram, [6](#)
- essHist (essHist-package), [2](#)
- essHist-package, [2](#)
- essHistogram, [4](#), [5](#), [8](#), [11](#), [12](#)
- essHistogram (Essential Histogram), [6](#)
- Generate Intervals, [8](#)
- genIntv, [5](#), [7](#), [11](#), [12](#)
- genIntv (Generate Intervals), [8](#)
- hist, [4](#), [7](#)
- Mixed normals, [9](#)
- mixnormal (Mixed normals), [9](#)
- msQuantile, [4–8](#)
- msQuantile (Multiscale Quantiles), [10](#)
- Multiscale Quantiles, [10](#)
- names, [11](#)
- Normal, [10](#)
- par, [4](#)
- paramExample (Mixed normals), [9](#)
- plot, [4](#)
- plot.histogram, [6](#)
- pmixnorm (Mixed normals), [9](#)
- pnorm, [9](#)
- quantile, [11](#)
- rmixnorm (Mixed normals), [9](#)
- rnorm, [9](#)
- suppressMessages, [5](#), [7](#), [11](#)
- title, [4](#), [6](#)