

Package ‘binr’

July 2, 2014

Title Cut Numeric Values Into Evenly Distributed Groups

Version 1.0

Author Sergei Izrailev

Maintainer Sergei Izrailev <sizrailev@collective.com>

Description This package provides algorithms for cutting numerical values exhibiting a potentially highly skewed distribution into evenly distributed groups (bins). This functionality can be applied for binning discrete values, such as counts, as well as for discretization of continuous values, for example, during generation of features used in machine learning algorithms.

URL <http://github.com/collectivemedia/binr>

Depends R (>= 3.0.3),

License Apache License (== 2.0)

Copyright Copyright (C) Collective, Inc. | file inst/COPYRIGHTS

LazyData true

R topics documented:

binr	2
bins	2
bins.greedy	5
bins.optimize	6
bins.quantiles	7
Index	9

`binr`*Cut Numeric Values Into Evenly Distributed Groups (bins).*

Description

Package `binr` (pronounced as "binner") provides algorithms for cutting numerical values exhibiting a potentially highly skewed distribution into evenly distributed groups (bins). This functionality can be applied for binning discrete values, such as counts, as well as for discretization of continuous values, for example, during generation of features used in machine learning algorithms.

Maintainer

Sergei Izrailev

Copyright

Copyright (C) Collective, Inc.

License

Apache License, Version 2.0, available at <http://www.apache.org/licenses/LICENSE-2.0>

URL

<http://github.com/collectivemedia/binr>

Installation from github

```
devtools::install_github("collectivemedia/binr")
```

Author(s)

Sergei Izrailev

See Also

[bins](#), [bins.quantiles](#), [bins.optimize](#), [bins.greedy](#)

`bins`*Cut Numeric Values Into Evenly Distributed Groups (Bins)*

Description

`bins` - Cuts points in vector `x` into evenly distributed groups (bins). `bins` takes 3 separate approaches to generating the cuts, picks the one resulting in the least mean square deviation from the ideal cut - $\text{length}(x) / \text{target.bins}$ points in each bin - and then merges small bins unless `exact.groups` is TRUE. The 3 approaches are:

1. Use quantiles, and increase the number of even cuts up to `max.breaks` until the number of groups reaches the desired number. See [bins.quantiles](#).
2. Start with a single bin with all the data in it and perform bin splits until either the desired number of bins is reached or there's no reduction in error (the latter is ignored if `exact.groups` is TRUE). See [bins.split](#).
3. Start with $\text{length}(\text{table}(x))$ bins, each containing exactly one distinct value and merge bins until the desired number of bins is reached. If `exact.groups` is FALSE, continue merging until there's no further reduction in error. See [bins.merge](#).

For each of these approaches, apply redistribution of points among existing bins until there's no further decrease in error. See `bins.move`.

`bins.getvals` - Extracts cut points from the object returned by `bins`. The cut points are always between the values in `x` and weighed such that the cut point splits the area under the line from (l_0, n_1) to (h_i, n_2) in half.

`bins.merr` - Partitioning the data into bins using splitting, merging and moving optimizes this error function, which is the mean squared error of point counts in the bins relative to the optimal number of points per bin.

Usage

```
bins(x, target.bins, max.breaks = NA, exact.groups = F, verbose = F,
     errthresh = 0.1, minpts = NA)
```

```
bins.getvals(lst, minpt = -Inf, maxpt = Inf)
```

```
bins.merr(binct, target.bins)
```

Arguments

<code>x</code>	Vector of numbers
<code>target.bins</code>	Number of groups desired; this is also the max number of groups.
<code>max.breaks</code>	Used for initial cut. If <code>exact.groups</code> is FALSE, bins are merged until there's no bins with fewer than $\text{length}(x) / \text{max.breaks}$ points. In <code>bins</code> , one of <code>max.breaks</code> and <code>minpts</code> must be supplied.
<code>exact.groups</code>	if TRUE, the result will have exactly the number of <code>target.bins</code> bins; if FALSE, the result may contain fewer than <code>target.bins</code> bins
<code>verbose</code>	Indicates verbose output.
<code>errthresh</code>	If the error is below the provided value, stops after the first rough estimate of the bins.
<code>minpts</code>	Minimum number of points in a bin. In <code>bins</code> , one of <code>max.breaks</code> and <code>minpts</code> must be supplied.
<code>lst</code>	The list returned by the <code>bins</code> function.
<code>minpt</code>	The value replacing the lower bound of the cut points.
<code>maxpt</code>	The value replacing the upper bound of the cut points.
<code>binct</code>	The number of points falling into the bins.

Details

The gains are computed using incremental analytical expressions derived for moving a value from one bin to the next, splitting a bin into two or merging two bins.

Value

A list containing the following items (not all of them may be present):

- `binlo` - The "low" value falling into the bin.
- `binhi` - The "high" value falling into the bin.
- `binct` - The number of points falling into the bin.
- `xtbl` - The result of a call to `table(x)`.
- `xval` - The sorted unique values of the data points `x`. Essentially, a numeric version of `names(xtbl)`.
- `changed` - Flag indicating whether the bins have been modified by the function.
- `err` - Mean square root error between the resulting counts and ideal bins.
- `imax` - For the move, merge and split operations, the index of the bin with the maximum gain.
- `iside` - For the move operation, the side of the move: 0 = left, 1 = right.
- `gain` - Error gain obtained as the result of the function call.

`bins.getvals` returns a vector of cut points extracted from the `lst` object.

See Also

[binr](#), [bins.greedy](#), [bins.quantiles](#) [bins.optimize](#)

Examples

```
## Not run:
# Seriously skewed x:
x <- floor(exp(rnorm(200000 * 1.3)))
cuts <- bins(x, target.bins = 10, minpts = 2000)
cuts$breaks <- bins.getvals(cuts)
cuts$binct
#   [0, 0]   [1, 1]   [2, 2]   [3, 3]   [4, 4]   [5, 5]   [6, 7]   [8, 10]
# 129868   66611   28039   13757   7595   4550   4623   2791
#   [11, 199]
# 2166

# Centered x:
x <- rep(c(1:10,20,31:40), c(rep(1, 10), 100, rep(1,10)))
cuts <- bins(x, target.bins = 3, minpts = 10)
cuts$binct
# [1, 10] [20, 20] [31, 40]
#      10      100      10

## End(Not run)
```

bins.greedy	<i>Greedy binning algorithm.</i>
-------------	----------------------------------

Description

`bins.greedy` - Wrapper around `bins.greedy.impl`. Goes over the sorted values of `x` left to right and fills the bins with the values until they are about the right size.

`bins.greedy.impl` - Implementation of a single-pass binning algorithm that examines sorted data left to right and builds bins of the target size. The `bins.greedy` wrapper around this function provides a less involved interface. This is not symmetric wrt direction: symmetric distributions may not have symmetric bins if there are multiple points with the same values. If a single value accounts for more than `thresh * binsz` points, it will be placed in a new bin.

Usage

```
bins.greedy(x, nbins, minpts = floor(0.5 * length(x)/nbins), thresh = 0.8,
  naive = FALSE)
```

```
bins.greedy.impl(xval, xtbl, xstp, binsz, nbins, thresh, verbose = F)
```

Arguments

<code>x</code>	Vector of numbers.
<code>nbins</code>	Target number of bins.
<code>minpts</code>	Minimum number of points in a bin. Only used if <code>naive = FALSE</code> .
<code>naive</code>	When <code>TRUE</code> , simply calls <code>bins.greedy.impl</code> with data derived from <code>x</code> . Otherwise, makes an extra step of marking the values that by themselves take a whole bin to force the algorithm to place these values in a bin separately.
<code>xval</code>	Sorted unique values of the data set <code>x</code> . This should be the numeric version of <code>names(xtbl)</code> .
<code>xtbl</code>	Result of a call to <code>table(x)</code> .
<code>xstp</code>	Stopping points; if <code>xstp[i] == TRUE</code> , the <i>i</i> -th value can't be merged to the (<i>i</i> -1)-th one. <code>xstp[1]</code> value is ignored.
<code>binsz</code>	Target bin size, i.e., the number of points falling into each bin; for example, <code>floor(length(x) / nbins)</code>
<code>thresh</code>	Threshold fraction of bin size for the greedy algorithm. Suppose there's <code>n < binsz</code> points in the current bin already. Also suppose that the next value <code>V</code> is represented by <code>m</code> points, and <code>m + n > binsz</code> . Then the algorithm will check if <code>m > thresh * binsz</code> , and if so, will place the value <code>V</code> into a new bin. If <code>m</code> is below the threshold, the points having value <code>V</code> are added to the current bin.
<code>verbose</code>	When <code>TRUE</code> , prints the number of points falling into the bins.

Value

A list with the following items:

- `binlo` - The "low" value falling into the bin.
- `binhi` - The "high" value falling into the bin.

- `binct` - The number of points falling into the bin.
- `xtbl` - The result of a call to `table(x)`.
- `xval` - The sorted unique values of the data points `x`. Essentially, a numeric version of `names(xtbl)`.

See Also

[binr](#), [bins](#), [bins.quantiles](#) [bins.optimize](#)

`bins.optimize`

Algorithms minimizing the binning error function `bins.merr`.

Description

`bins.move` - Compute the best move of a value from one bin to its neighbor
`bins.split` - Split a bin into two bins optimally.
`bins.merge` - Merges the two bins yielding the largest gain in error reduction.
`bins.move.iter` - Apply `bins.move` until there's no change. Can only reduce the error.
`bins.split.iter` Iterate to repeatedly apply `bins.split`.
`bins.merge.iter` Iterate to repeatedly apply `bins.merge`.

Usage

```
bins.move(xval, xtbl, binlo, binhi, binct, target.bins, verbose = F)

bins.split(xval, xtbl, binlo, binhi, binct, target.bins, force = F,
  verbose = F)

bins.merge(xval, xtbl, binlo, binhi, binct, target.bins, force = F,
  verbose = F)

bins.move.iter(lst, target.bins, verbose = F)

bins.split.iter(lst, target.bins, exact.groups = F, verbose = F)

bins.merge.iter(lst, target.bins, exact.groups = F, verbose = F)
```

Arguments

<code>xval</code>	Sorted unique values of the data set <code>x</code> . This should be the numeric version of <code>names(xtbl)</code> .
<code>xtbl</code>	Result of a call to <code>table(x)</code> .
<code>binlo</code>	The "low" value falling into the bin.
<code>binhi</code>	The "high" value falling into the bin.
<code>binct</code>	The number of points falling into the bin.
<code>target.bins</code>	Number of bins desired; this is also the max number of bins.
<code>force</code>	When TRUE, splits or merges bins regardless of whether the best gain is positive.

<code>verbose</code>	When <code>TRUE</code> , prints resulting <code>binct</code> .
<code>lst</code>	List containing <code>xval</code> , <code>xtbl</code> , <code>binlo</code> , <code>binhi</code> , <code>binct</code> .
<code>exact.groups</code>	If <code>FALSE</code> , run until either the <code>target.bins</code> is reached or there's no more splits or merges that reduce the error. Otherwise (<code>TRUE</code>), run until the <code>target.bins</code> is reached, even if that increases the error.

Value

A list containing the following items (not all of them may be present):

- `binlo` - The "low" value falling into the bin.
- `binhi` - The "high" value falling into the bin.
- `binct` - The number of points falling into the bin.
- `xtbl` - The result of a call to `table(x)`.
- `xval` - The sorted unique values of the data points `x`. Essentially, a numeric version of `names(xtbl)`.
- `changed` - Flag indicating whether the bins have been modified by the function.
- `err` - Mean square root error between the resulting counts and ideal bins.
- `imax` - For the move, merge and split operations, the index of the bin with the maximum gain.
- `iside` - For the move operation, the side of the move: 0 = left, 1 = right.
- `gain` - Error gain obtained as the result of the function call.

See Also

[bins](#), [binr](#), [bins.greedy](#), [bins.quantiles](#)

<code>bins.quantiles</code>	<i>Quantile-based binning</i>
-----------------------------	-------------------------------

Description

Cuts the data set `x` into roughly equal groups using quantiles.

Usage

```
bins.quantiles(x, target.bins, max.breaks, verbose = FALSE)
```

Arguments

<code>x</code>	A numeric vector to be cut in bins.
<code>target.bins</code>	Target number of bins, which may not be reached if the number of unique values is smaller than the specified value.
<code>max.breaks</code>	Maximum number of quantiles; must be at least as large as <code>target.bins</code> .
<code>verbose</code>	Indicates verbose output.

Details

Because the number of unique values may be smaller than `target.bins`, the function gradually increases the number of quantiles up to `max.breaks` or until the `target.bins` number of bins is reached.

See Also

[binr](#), [bins](#), [bins.greedy](#), [bins.optimize](#)

Index

*Topic **64-bit**

binr, 2

*Topic **bigint**

binr, 2

*Topic **csv**

binr, 2

*Topic **delimited**

binr, 2

*Topic **file**

binr, 2

*Topic **integer64**

binr, 2

*Topic **read.csv**

binr, 2

binr, 2, 4, 6–8

binr-package (*binr*), 2

bins, 2, 2, 6–8

bins.greedy, 2, 4, 5, 7, 8

bins.merge, 3

bins.merge (*bins.optimize*), 6

bins.move (*bins.optimize*), 6

bins.optimize, 2, 4, 6, 6, 8

bins.quantiles, 2–4, 6, 7, 7

bins.split, 3

bins.split (*bins.optimize*), 6