

Package ‘salso’

August 13, 2021

Type Package

Title Search Algorithms and Loss Functions for Bayesian Clustering

Version 0.2.23

Description The SALSO algorithm is an efficient randomized greedy search method to find a point estimate for a random partition based on a loss function and posterior Monte Carlo samples. The algorithm is implemented for many loss functions, including the Binder loss and a generalization of the variation of information loss, both of which allow for unequal weights on the two types of clustering mistakes. Efficient implementations are also provided for Monte Carlo estimation of the posterior expected loss of a given clustering estimate. See Dahl, Johnson, Müller (2021) <[arXiv:2105.04451](https://arxiv.org/abs/2105.04451)>.

License MIT + file LICENSE | Apache License 2.0

Depends R (>= 4.0.0)

SystemRequirements Cargo (>= 1.51.0) for installation from sources:
see INSTALL file

Encoding UTF-8

LazyData TRUE

RoxygenNote 7.1.1

NeedsCompilation yes

Author David B. Dahl [aut, cre] (<<https://orcid.org/0000-0002-8173-1547>>),
Devin J. Johnson [aut] (<<https://orcid.org/0000-0003-2619-6649>>),
Peter Müller [aut]

Maintainer David B. Dahl <dahl@stat.byu.edu>

Repository CRAN

Date/Publication 2021-08-13 04:40:02 UTC

R topics documented:

bell	2
dlso	3
enumerate.partitions	4
enumerate.permutations	4

iris.clusterings	5
partition.loss	6
plot.salso.summary	10
psm	11
salso	12
summary.salso.estimate	14

Index	16
--------------	-----------

bell	<i>Compute the Bell Number</i>
------	--------------------------------

Description

These functions compute the Bell number (the number of partitions of a set) or its natural logarithm.

Usage

```
bell(nItems)
```

```
lbell(nItems)
```

Arguments

nItems The size of the set 1, 2, ..., n .

Value

A numeric vector of length one giving the Bell number or its natural logarithm.

Examples

```
bell(12)
lbell(300)
all.equal( bell(5), exp(lbell(5)) )
```

Description

This function provides a partition to summarize a partition distribution using the draws-based latent structure optimization (DLSO) method, which is also known as the least-squares clustering method (Dahl 2006). The method seeks to minimize an estimation criterion by picking the minimizer among the partitions supplied. The implementation currently supports the minimization of several partition estimation criteria. For details on these criteria, see [partition.loss](#).

Usage

```
dlsO(partitions, loss = VI(), x = NULL)
```

Arguments

partitions	An integer matrix of cluster labels with n columns, where each row is a partition of n items given as cluster labels. Two items are in the same subset (i.e., cluster) if their labels are equal. Or, a vector of length n which will be converted to a 1-by- n matrix.
loss	The loss function to use, as indicated by "binder", "omARI", "VI", "NVI", "ID", "NID", or the result of calling a function with these names. Also supported are "binder.psm", "VI.lb", "omARI.approx", or the result of calling a function with these names, in which case x above can optionally be a pairwise similarity matrix, i.e., n -by- n symmetric matrix whose (i, j) element gives the (estimated) probability that items i and j are in the same subset (i.e., cluster) of a partition (i.e., clustering).
x	A B -by- n matrix, where each of the B rows represents a clustering of n items using cluster labels. For the b th clustering, items i and j are in the same cluster if $x[b, i] == x[b, j]$.

Value

An integer vector giving the estimated partition, encoded using cluster labels.

References

D. B. Dahl (2006), Model-Based Clustering for Expression Data via a Dirichlet Process Mixture Model, in *Bayesian Inference for Gene Expression and Proteomics*, Kim-Anh Do, Peter Müller, Marina Vannucci (Eds.), Cambridge University Press.

See Also

[partition.loss](#), [psm](#), [summary.salso.estimate](#), [salso](#)

Examples

```
dls(iris.clusterings, loss=VI())
dls(iris.clusterings, loss=binder())

# Compute expected loss using all draws, but pick the best among the first 10.
dls(iris.clusterings[1:10,], loss=VI(), x=iris.clusterings)
```

enumerate.partitions *Enumerate Partitions of a Set*

Description

This function produces a matrix whose rows provide all possible partitions of the set $1, 2, \dots, n$. These partitions are provided as cluster labels, where two items are in the same subset (i.e., cluster) if their labels are equal.

Usage

```
enumerate.partitions(nItems)
```

Arguments

nItems The size of the set $1, 2, \dots, n$, i.e., n .

Value

A matrix of integers, where each row is a partition encoded as a vector of cluster labels.

Examples

```
enumerate.partitions(5)
```

enumerate.permutations
 Enumerate Permutations of Items

Description

This function produces a matrix whose rows provide all possible permutations of the set $1, 2, \dots, n$.

Usage

```
enumerate.permutations(nItems)
```

Arguments

`nItems` The size of the set $1, 2, \dots, n$, i.e., n .

Value

A matrix of integers, where each row is a permutation.

Examples

```
enumerate.permutations(5)
```

`iris.clusterings` *Clusterings of the Iris Data*

Description

Randomly generated clusterings for the `iris` dataset.

Usage

```
iris.clusterings
```

Format

A 1000-by-150 matrix of 1000 randomly generated clusterings of the 150 observations in the `iris` dataset.

Source

Unknown.

See Also

[iris](#)

 partition.loss

Compute Partition Loss or the Expectation of Partition Loss

Description

Given two partitions π^* and π , these functions compute the specified loss when using π^* to estimate π . Smaller loss values indicate higher concordance between partitions. These functions also compute a Monte Carlo estimate of the expectation for the specified loss based on samples or a pairwise similarity matrix. This function also supports computing approximations to the expectation of several losses. Supported criteria are described below. Some criteria only require the pairwise similarity matrix (as computed, for example, by [psm](#)) whereas others require samples from a partition distribution. For those criteria that only need the pairwise similarity matrix, posterior samples can still be provided in the `x` argument and the pairwise similarity matrix will automatically be computed as needed.

Usage

```
partition.loss(partitions, x, loss = VI())
```

```
binder(partitions, x, a = 1)
```

```
RI(partition1, partition2)
```

```
omARI(partitions, x)
```

```
omARI.approx(partitions, x)
```

```
ARI(partition1, partition2)
```

```
VI(partitions, x, a = 1)
```

```
VI.lb(partitions, x)
```

```
NVI(partitions, x)
```

```
ID(partitions, x)
```

```
NID(partitions, x)
```

Arguments

partitions	An integer matrix of cluster labels with n columns, where each row is a partition of n items given as cluster labels. Two items are in the same subset (i.e., cluster) if their labels are equal. Or, a vector of length n which will be converted to a 1-by- n matrix.
------------	---

x	A B -by- n matrix, where each of the B rows represents a clustering of n items using cluster labels. For the b th clustering, items i and j are in the same cluster if $x[b, i] == x[b, j]$.
loss	The loss function to use, as indicated by "binder", "omARI", "VI", "NVI", "ID", "NID", or the result of calling a function with these names. Also supported are "binder.psm", "VI.lb", "omARI.approx", or the result of calling a function with these names, in which case x above can optionally be a pairwise similarity matrix, i.e., n -by- n symmetric matrix whose (i, j) element gives the (estimated) probability that items i and j are in the same subset (i.e., cluster) of a partition (i.e., clustering).
a	(Only used for Binder and VI loss) Without loss of generality, the cost under Binder loss of placing two items in the same cluster when in truth they belong to separate clusters is fixed at one. The argument a is either: i. a nonnegative scalar giving the cost of the complementary mistake, i.e., placing two items in separate clusters when in truth they belong to the same cluster, or ii. a list containing the desired number of clusters ("nClusters") and an upper bound ("upper") when searching for " a " that yields the desired number of clusters. In the second case, the desired number of clusters may not be achievable without modifying <code>maxSize</code> in the <code>salso</code> function. To increase the probability of hitting exactly the desired number of clusters, the <code>nRuns</code> in the <code>salso</code> function may need to be increased.
partition1	An integer vector of cluster labels for n items. Two items are in the same subset (i.e., cluster) if their labels are equal.
partition2	An integer vector of cluster labels having the same length as <code>partition1</code> .

Details

The partition estimation criterion can be specified using the `loss` argument, which is either a string or a result of calling the associated functions. These losses are described below:

"binder" Binder. Whereas high values of the Rand index R between π^* and π correspond to high concordance between the partitions, the N-invariant Binder loss L for a partition π^* in estimating π is $L = (1 - R) * (n - 1)/n$, meaning that low values correspond to high concordance between the partitions. This package reports the N-invariant Binder loss. The original Binder loss equals the N-invariant Binder loss multiplied by $n^2/2$. Only the pairwise similarity matrix is required for this loss, but samples can be provided. As originally discussed by Binder (1978), two mistakes are possible: 1. Placing two items in separate clusters when in truth they belong to the same cluster, and 2. Placing two items in the same cluster when in truth they belong to separate clusters. Without loss of generality, the cost of the second mistake is fixed at one. The default cost of the first mistake is also one, but can be specified with the argument a . For a discussion of general weights, see Dahl, Johnson, and Müller (2021). For a discussion of the equal weights case, see also Dahl (2006), Lau and Green (2007), Dahl and Newton (2007), Fritsch and Ickstadt (2009), and Wade and Ghahramani (2018).

"omARI" One Minus Adjusted Rand Index. Computes the expectation of one minus the adjusted Rand index (Hubert and Arabie, 1985). Whereas high values of the adjusted Rand index between π^* and π correspond to high concordance between the partitions, the loss associated with the adjusted Rand index for a partition π^* in estimating π is one minus the adjusted Rand index between the partitions, meaning that low values correspond to high concordance

between the partitions. Samples from a partition distribution are required for this loss. See Fritsch and Ickstadt (2009).

"omARI.approx" Approximation of One Minus Adjusted Rand Index. Computes the first-order approximation of the expectation of one minus the adjusted Rand index. The adjusted Rand index involves a ratio and the first-order approximation of the expectation is based on $E(X/Y) \approx E(X)/E(Y)$. Only the pairwise similarity matrix is required for this criterion, but samples can be provided. See Fritsch and Ickstadt (2009).

"VI" Variation of Information. Computes the expectations of the (generalized) variation of information loss. Samples from a partition distribution are required for this loss. See Meilă (2007), Wade and Ghahramani (2018), and Rastelli and Friel (2018). The original variation of information of Meilă (2007) has been extended to the generalized variation of information of Dahl, Johnson, and Müller (2021) to allow for unequal weighting of two possible mistakes: 1. Placing two items in separate clusters when in truth they belong to the same cluster, and 2. Placing two items in the same cluster when in truth they belong to separate clusters. Without loss of generality, the weight for the second mistake is fixed at one. The default weight of the first mistake is also one, but can be specified with the argument *a*. See Dahl, Johnson, Müller (2021).

"VI.lb" Lower Bound of the Variation of Information. Computes the lower bound of the expectation of the variation of information loss, where the lower bound is obtained by Jensen's inequality. Only the pairwise similarity matrix is required for this criterion, but samples can be provided. See Wade and Ghahramani (2018).

"NVI" Normalized Variation of Information. Computes the expectation of the normalized variation of information loss. Samples from a partition distribution are required for this loss. See Vinh, Epps, and Bailey (2010) and Rastelli and Friel (2018).

"ID" Information Distance. Computes the expectation of the information distance (D_{max}) loss. Samples from a partition distribution are required for this loss. See Vinh, Epps, and Bailey (2010).

"NID" Normalized Information Distance. Computes the expectation of the normalized information distance loss. Samples from a partition distribution are required for this loss. See Vinh, Epps, and Bailey (2010) and Rastelli and Friel (2018).

The functions [RI](#) and [ARI](#) are convenience functions. Note that:

- $binder(p1, p2) = (1 - RI(p1, p2)) * (n - 1) / n$
- $omARI(p1, p2) = 1 - ARI(p1, p2)$

Value

A numeric vector.

References

- W. M. Rand (1971), Objective Criteria for the Evaluation of Clustering Methods. *Journal of the American Statistical Association*, **66**: 846–850.
- D. A. Binder (1978), Bayesian cluster analysis, *Biometrika* **65**, 31-38.
- L. Hubert and P. Arabie (1985), Comparing Partitions. *Journal of Classification*, **2**, 193–218.

- D. B. Dahl (2006), Model-Based Clustering for Expression Data via a Dirichlet Process Mixture Model, in *Bayesian Inference for Gene Expression and Proteomics*, Kim-Anh Do, Peter Müller, Marina Vannucci (Eds.), Cambridge University Press.
- J. W. Lau and P. J. Green (2007), Bayesian model based clustering procedures, *Journal of Computational and Graphical Statistics* **16**, 526-558.
- M. Meilă (2007), Comparing Clusterings - an Information Based Distance. *Journal of Multivariate Analysis*, **98**: 873–895.
- D. B. Dahl and M. A. Newton (2007), Multiple Hypothesis Testing by Clustering Treatment Effects, *Journal of the American Statistical Association*, **102**, 517-526.
- A. Fritsch and K. Ickstadt (2009), An improved criterion for clustering based on the posterior similarity matrix, *Bayesian Analysis*, **4**, 367-391.
- N. X. Vinh, J. Epps, and J. Bailey (2010), Information Theoretic Measures for Clusterings Comparison: Variants, Properties, Normalization and Correction for Chance, *Journal of Machine Learning Research*, **11**, 2837-2854.
- S. Wade and Z. Ghahramani (2018), Bayesian cluster analysis: Point estimation and credible balls. *Bayesian Analysis*, **13:2**, 559-626.
- R. Rastelli and N. Friel (2018), Optimal Bayesian estimators for latent variable cluster models. *Statistics and Computing*, **28**, 1169-1186.
- D. B. Dahl, D. J. Johnson, and P. Müller (2021), Search Algorithms and Loss Functions for Bayesian Clustering, <arXiv:2105.04451>.

See Also

[psm](#), [salso](#), [dlso](#)

Examples

```
# For examples, use 'nCores=1' per CRAN rules, but in practice omit this.
partitions <- iris.clusterings[1:5,]

all.equal(partition.loss(partitions, partitions, loss=binder(a=2)),
          binder(partitions, partitions, a=2))
all.equal(partition.loss(partitions, partitions, loss=omARI()),
          omARI(partitions, partitions))
all.equal(partition.loss(partitions, partitions, loss=VI(a=0.8)),
          VI(partitions, partitions, a=0.8))

p1 <- iris.clusterings[1,]
p2 <- iris.clusterings[2,]

VI(p1, p2, a=1.0)
all.equal(binder(p1, p2), ( 1 - RI(p1, p2) ) * (length(p1)-1) / length(p1))
all.equal(omARI(p1, p2), 1 - ARI(p1, p2))
```

plot.salso.summary	<i>Heatmap, Multidimensional Scaling, Pairs, and Dendrogram Plotting for Partition Estimation</i>
--------------------	---

Description

This function produces one of four plots: 1. "heatmap": A heatmap showing the pairwise allocation probabilities that items are clustered. 2. "mds": A scatter plot using classical multidimensional scaling (also known as principal coordinates analysis) with the exemplar (i.e., the most representative observation) of each cluster emphasized. 3. "pairs": Pairs plots of all the variables with the exemplar (i.e., the most representative observation) of each cluster emphasized. 4. "dendrogram": A dendrogram based on expected partition loss showing the relationships among clusters when merging pairs of clusters such that the increase in the expectation of the posterior loss is minimized.

Usage

```
## S3 method for class 'salso.summary'
plot(
  x,
  type = c("heatmap", "mds", "pairs", "dendrogram")[1],
  data = NULL,
  showLabels = TRUE,
  showIDs = length(x$estimate) <= 50,
  ...
)
```

Arguments

x	An object returned by <code>summary(y)</code> , where <code>y</code> itself is returned by the salso function.
type	A string equal to "heatmap", "mds", "pairs", or "dendrogram".
data	The data from which the partition estimation was obtained. This is required when <code>type='pairs'</code> and ignored otherwise.
showLabels	Should the cluster labels be shown in the plot when <code>type="heatmap"</code> ?
showIDs	Should the ID of the items be shown in the plot?
...	Arguments to be passed to methods, such as graphical parameters (see par).

Value

NULL, invisibly.

See Also

[salso](#), [summary.salso.estimate](#), [cmdscale](#).

Examples

```
# For examples, use 'nCores=1' per CRAN rules, but in practice omit this.
draws <- iris.clusterings
est <- salso(draws, nCores=1)
summ <- summary(est)
plot(summ, type="heatmap")
plot(summ, type="mds")
plot(summ, type="pairs", data=iris)
plot(summ, type="dendrogram")
```

 psm

Compute an Adjacency or Pairwise Similarity Matrix

Description

If only one sample is provided, this function computes an adjacency matrix, i.e., a binary matrix whose (i, j) element is one if and only if elements i and j in the partition have the same cluster label. If multiple samples are provided (as rows of the x matrix), this function computes the n -by- n matrix whose (i, j) element gives the relative frequency (i.e., estimated probability) that items i and j are in the same subset (i.e., cluster). This is the mean of the adjacency matrices of the provided samples.

Usage

```
psm(x, nCores = 0)
```

Arguments

x	A B -by- n matrix, where each of the B rows represents a clustering of n items using cluster labels. For the b th clustering, items i and j are in the same cluster if $x[b, i] == x[b, j]$.
nCores	The number of CPU cores to use, i.e., the number of simultaneous runs at any given time. A value of zero indicates to use all cores on the system.

Value

A n -by- n symmetric matrix whose (i, j) element gives the relative frequency that that items i and j are in the same subset (i.e., cluster).

Examples

```
# For examples, use 'nCores=1' per CRAN rules, but in practice omit this.

partition <- iris.clusterings[1,]
psm(partition, nCores=1)

dim(iris.clusterings)
```

```
probs <- psm(iris.clusterings, nCores=1)
dim(probs)
probs[1:6, 1:6]
```

salso

SALSO Greedy Search

Description

This function provides a partition to summarize a partition distribution using the SALSO greedy search method (Dahl, Johnson, and Müller, 2021). The implementation currently supports the minimization of several partition estimation criteria. For details on these criteria, see [partition.loss](#).

Usage

```
salso(
  x,
  loss = VI(),
  maxNClusters = 0,
  nRuns = 16,
  maxZealousAttempts = 10,
  probSequentialAllocation = 0.5,
  nCores = 0,
  ...
)
```

Arguments

- | | |
|---------------------------|---|
| <code>x</code> | A B -by- n matrix, where each of the B rows represents a clustering of n items using cluster labels. For the b th clustering, items i and j are in the same cluster if <code>x[b, i] == x[b, j]</code> . |
| <code>loss</code> | The loss function to use, as indicated by "binder", "omARI", "VI", "NVI", "ID", "NID", or the result of calling a function with these names. Also supported are "binder.psm", "VI.lb", "omARI.approx", or the result of calling a function with these names, in which case <code>x</code> above can optionally be a pairwise similarity matrix, i.e., n -by- n symmetric matrix whose (i, j) element gives the (estimated) probability that items i and j are in the same subset (i.e., cluster) of a partition (i.e., clustering). The loss functions "binder.psm", "VI.lb", and "omARI.approx" are generally not recommended and the current implementation requires that <code>maxZealousAttempts = 0</code> and <code>probSequentialAllocation = 1.0</code> . |
| <code>maxNClusters</code> | The maximum number of clusters that can be considered by the optimization algorithm, which has important implications for the interpretability of the resulting clustering and can greatly influence the RAM needed for the optimization algorithm. If the supplied value is zero and <code>x</code> is a matrix of clusterings, the |

optimization is constrained by the maximum number of clusters among the clusterings in x . If the supplied value is zero and x is a pairwise similarity matrix, there is no constraint.

<code>nRuns</code>	The number of runs to try, although the actual number may differ for the following reasons: 1. The actual number is a multiple of the number of cores specified by the <code>nCores</code> argument, and 2. The search is curtailed when the seconds threshold is exceeded.
<code>maxZealousAttempts</code>	The maximum number of attempts for zealous updates, in which entire clusters are destroyed and items are sequentially reallocated. While zealous updates may be helpful in optimization, they also take more CPU time which might be better used trying additional runs.
<code>probSequentialAllocation</code>	For the initial allocation, the probability of sequential allocation instead of using <code>sample(1:K, ncol(x), TRUE)</code> , where K is set according to the <code>maxNClusters</code> argument.
<code>nCores</code>	The number of CPU cores to use, i.e., the number of simultaneous runs at any given time. A value of zero indicates to use all cores on the system.
<code>...</code>	Extra arguments not intended for the end user, including: 1. <code>seconds</code> : Instead of performing all the requested number of runs, curtail the search after the specified expected number of seconds. Note that the function will finish earlier if all the requested runs are completed. The specified seconds does not account for the overhead involved in starting the search and returning results. 2. <code>maxScans</code> : The maximum number of full reallocation scans. The actual number of scans may be less than <code>maxScans</code> since the method stops if the result does not change between scans, and 3. <code>probSingletonsInitialization</code> : When doing a sequential allocation to obtain the initial allocation, the probability of placing the first <code>maxNClusters</code> randomly-selected items in singletons subsets.

Details

The initial version of the SALSO method was presented at the workshop "Bayesian Nonparametric Inference: Dependence Structures and their Applications" in Oaxaca, Mexico on December 6, 2017. See <https://www.birs.ca/events/2017/5-day-workshops/17w5060/schedule>.

Value

An integer vector giving the estimated partition, encoded using cluster labels.

References

D. B. Dahl, D. J. Johnson, and P. Müller (2021), Search Algorithms and Loss Functions for Bayesian Clustering, [arXiv:2105.04451](https://arxiv.org/abs/2105.04451).

See Also

[partition.loss](#), [psm](#), [summary.salso.estimate](#), [dlso](#)

Examples

```
# For examples, use 'nCores=1' per CRAN rules, but in practice omit this.

draws <- iris.clusterings
salso(draws, loss=VI(), nRuns=1, nCores=1)
salso(draws, loss=VI(a=0.7), nRuns=1, nCores=1)
salso(draws, loss=binder(), nRuns=1, nCores=1)
salso(iris.clusterings, binder(a=list(nClusters=3, upper=5)), nRuns=4, nCores=1)
```

```
summary.salso.estimate
```

Summary of Partitions Estimated Using Posterior Expected Loss

Description

Assessing the quality of clusters in a partition estimate is added by this function. The result can then be plotted with `plot.salso.summary`. The current implementation of the calculation of these summaries is not terribly efficient and may be improved in the future.

Usage

```
## S3 method for class 'salso.estimate'
summary(object, alternative, orderingMethod = 1, ...)
```

Arguments

<code>object</code>	An object returned by the <code>salso</code> function.
<code>alternative</code>	An optional argument specifying an alternative clustering to use instead of that provided by <code>object</code> . Use this feature to obtain numerical and graphical summaries of a clustering estimate from other procedures. This clustering must be provided in canonical form: cluster labels as integers starting at 1 for the first observation and incrementing by one for each new label.
<code>orderingMethod</code>	An integer giving method to use to order the observations for a heatmap plot. Currently values 1 or 2 are supported.
<code>...</code>	Currently ignored.

Value

A list containing the estimate, the pairwise similarity matrix, the mean pairwise similarity matrix, the score and mean pairwise similarity for each observation, exemplar observation for each cluster, a dendrogram object, a vector for ordering observations in the heatmap plot, the size of each cluster, and the number of clusters.

Examples

```
# For examples, use 'nCores=1' per CRAN rules, but in practice omit this.
draws <- iris.clusterings
est <- salso(draws, nCores=1)
summ <- summary(est)
plot(summ, type="heatmap")
plot(summ, type="mds")
plot(summ, type="pairs", data=iris)
plot(summ, type="dendrogram")
```

Index

* datasets

iris.clusterings, 5

ARI, 8

ARI (partition.loss), 6

bell, 2

binder (partition.loss), 6

cmdscafe, 10

dlso, 3, 9, 13

enumerate.partitions, 4

enumerate.permutations, 4

ID (partition.loss), 6

iris, 5

iris.clusterings, 5

lbelle (bell), 2

NID (partition.loss), 6

NVI (partition.loss), 6

omARI (partition.loss), 6

par, 10

partition.loss, 3, 6, 12, 13

plot.salso.summary, 10, 14

psm, 3, 6, 9, 11, 13

RI, 8

RI (partition.loss), 6

salso, 3, 7, 9, 10, 12, 14

summary.salso.estimate, 3, 10, 13, 14

VI (partition.loss), 6