

Package ‘mSigHdp’

June 6, 2020

Title Mutational signature extraction using hdp (Hierarchical Dirichlet Process)

Version 0.0.0.9015

Description Calls hdp for mutational signature analysis, with performance issues in hdp::stirling() corrected.

License GPL-3

Encoding UTF-8

LazyData true

Language en-US

biocViews

Imports hdp (>= 0.1.5.0014),
SynSigGen

Roxygen list(markdown = TRUE)

Depends R (>= 3.5)

RoxygenNote 7.1.0

Remotes github::steverozen/hdp,
github::steverozen/SynSigGen,
github::WuyangFF95/SynSigEval

Suggests testthat,
ICAMS,
utils,
SynSigEval

R topics documented:

AnalyzeAndPlotretval	2
ChainsDiagnosticPlot	2
CombinePosteriorChains	3
ExtendIterationAndPosterior	4
MultipleSetupAndPosterior	6
PlotExposure	7
PlotExposureByRange	8
PrepInit	8
RunHdpParallel	9
SetupAndPosterior	10
SortExp	12
Index	13

AnalyzeAndPlotretval	<i>Evaluate and plot retval from CombinePosteriorChains</i>
----------------------	---

Description

Evaluate and plot retval from CombinePosteriorChains

Usage

```
AnalyzeAndPlotretval(
  retval,
  out.dir = NULL,
  ground.truth.sig = NULL,
  ground.truth.exp = NULL,
  verbose = TRUE,
  overwrite = TRUE
)
```

Arguments

retval	the output from function CombinePosteriorChains
out.dir	Directory that will be created for the output; if overwrite is FALSE then abort if out.dir already exists.
ground.truth.sig	Optional. Either a string with the path to file with ground truth signatures or and ICAMS catalog with the ground truth signatures. These are the signatures used to construct the ground truth spectra.
ground.truth.exp	Optional. Ground truth exposure matrix or path to file with ground truth exposures. If NULL skip checks that need this information.
verbose	If TRUE then message progress information.
overwrite	If TRUE overwrite out.dir if it exists, otherwise raise an error.

ChainsDiagnosticPlot	<i>Diagnostic plot for hdp multi sample chains (output from CombinePosteriorChains)</i>
----------------------	---

Description

Diagnostic plot for hdp multi sample chains (output from CombinePosteriorChains)

Usage

```
ChainsDiagnosticPlot(retval, out.dir, verbose)
```

Arguments

retval	output from CombinePosteriorChains. A list with the following elements: signature The extracted signature profiles as a matrix; rows are mutation types, columns are samples (e.g. tumors). exposure The inferred exposures as a matrix of mutation counts; rows are signatures, columns are samples (e.g. tumors). multi.chains A <code>hdpSampleMulti-class</code> object. This object has the method <code>chains</code> which returns a list of <code>hdpSampleChain-class</code> objects. Each of these sample chains objects has a method <code>final_hdpState</code> (actually the methods seems to be just <code>hdp</code>) that returns the <code>hdpState</code> from which it was generated.
out.dir	Directory that will be created for the output; if overwrite is FALSE then abort if out.dir already exists.
verbose	If TRUE then message progress information.

CombinePosteriorChains

Extract components and exposures from multiple posterior sample chains

Description

Extract components and exposures from multiple posterior sample chains

Usage

```
CombinePosteriorChains(
  clean.chlist,
  input.catalog,
  multi.types,
  verbose = TRUE,
  cos.merge = 0.9,
  min.sample = 1
)
```

Arguments

clean.chlist	It collects the output of multiple independent <code>hdp_posterior</code> calls.
input.catalog	Input spectra catalog as a matrix or in ICAMS format.
multi.types	A logical scalar or a character vector. If FALSE, <code>hdp</code> will regard all input spectra as one tumor type. If TRUE, <code>hdp</code> will infer tumor types based on the string before ":" in their names. e.g. tumor type for "SA.Syn.Ovary-AdenoCA::S.500" would be "SA.Syn.Ovary-AdenoCA" If <code>multi.types</code> is a character vector, then it should be of the same length as the number of columns in <code>input.catalog</code> , and each value is the name of the tumor type of the corresponding column in <code>input.catalog</code> , e.g. <code>c("SA.Syn.Ovary-AdenoCA", "SA.Syn.Ovary-AdenoCA")</code>
verbose	If TRUE then message progress information.

<code>cos.merge</code>	The cosine similarity threshold for merging raw clusters from the posterior sampling chains into "components" i.e. signatures; passed to hdp_extract_components .
<code>min.sample</code>	A "component" (i.e. signature) must have at least this many samples; passed to hdp_extract_components .

Value

Invisibly, a list with the following elements:

signature The extracted signature profiles as a matrix; rows are mutation types, columns are samples (e.g. tumors).

exposure The inferred exposures as a matrix of mutation counts; rows are signatures, columns are samples (e.g. tumors).

multi.chains A [hdpSampleMulti-class](#) object. This object has the method [chains](#) which returns a list of [hdpSampleChain-class](#) objects. Each of these sample chains objects has a method [final_hdpState](#) (actually the methods seems to be just hdp) that returns the hdpState from which it was generated.

ExtendIterationAndPosterior

Run hdp extraction and attribution on a spectra catalog file This repeats what Nicola do in her thesis. It starts four independent initial chains with `post.burnin` iterations, then pick up from the end of each initial chain, started another `num.posterior` MCMC chains for another `post.burnin` iterations and then collected `post.n` posterior samples at intervals of `post.space` iterations. In total, this collects 4 times `num.posterior` times `post.n` posterior samples from 4 times y separate chains.

Description

Run hdp extraction and attribution on a spectra catalog file This repeats what Nicola do in her thesis. It starts four independent initial chains with `post.burnin` iterations, then pick up from the end of each initial chain, started another `num.posterior` MCMC chains for another `post.burnin` iterations and then collected `post.n` posterior samples at intervals of `post.space` iterations. In total, this collects 4 times `num.posterior` times `post.n` posterior samples from 4 times y separate chains.

Usage

```
ExtendIterationAndPosterior(
  input.catalog,
  CPU.cores = 1,
  seedNumber = 1,
  K.guess,
  multi.types = FALSE,
  verbose = TRUE,
  num.posterior = 4,
  post.burnin = 4000,
  post.n = 50,
  post.space = 50,
  post.cpiter = 3,
```

```

    post.verbosity = 0,
    cos.merge = 0.9,
    min.sample = 1,
    checkpoint.aft.post = NULL
  )

```

Arguments

<code>input.catalog</code>	Input spectra catalog as a matrix or in ICAMS format.
<code>CPU.cores</code>	Number of CPUs to use in running hdp_posterior ; this is used to parallelize running the posterior sampling chains, so there is no point in making this larger than <code>num.posterior</code> .
<code>seedNumber</code>	An integer that is used to generate separate random seeds for each call to dp_activate , and each call of hdp_posterior ; please see the code on how this is done. But repeated calls with same value of <code>seedNumber</code> and other inputs should produce the same results.
<code>K.guess</code>	Suggested initial value of the number of signatures, passed to dp_activate as <code>initcc</code> .
<code>multi.types</code>	A logical scalar or a character vector. If <code>FALSE</code> , <code>hdp</code> will regard all input spectra as one tumor type. If <code>TRUE</code> , <code>hdp</code> will infer tumor types based on the string before ":" in their names. e.g. tumor type for "SA.Syn.Ovary-AdenoCA::S.500" would be "SA.Syn.Ovary-AdenoCA" If <code>multi.types</code> is a character vector, then it should be of the same length as the number of columns in <code>input.catalog</code> , and each value is the name of the tumor type of the corresponding column in <code>input.catalog</code> , e.g. <code>c("SA.Syn.Ovary-AdenoCA", "SA.Syn.O")</code>
<code>verbose</code>	If <code>TRUE</code> then message progress information.
<code>num.posterior</code>	Number of posterior sampling chains; can set to 1 for testing.
<code>post.burnin</code>	Pass to hdp_posterior <code>burnin</code> .
<code>post.n</code>	Pass to hdp_posterior <code>n</code> .
<code>post.space</code>	Pass to hdp_posterior <code>space</code> .
<code>post.cpiter</code>	Pass to hdp_posterior <code>cpiter</code> .
<code>post.verbosity</code>	Pass to hdp_posterior <code>verbosity</code> .
<code>cos.merge</code>	The cosine similarity threshold for merging raw clusters from the posterior sampling chains into "components" i.e. signatures; passed to hdp_extract_components .
<code>min.sample</code>	A "component" (i.e. signature) must have at least this many samples; passed to hdp_extract_components .
<code>checkpoint.aft.post</code>	If non-NULL, a file path to checkpoint the list of values returned from the calls to hdp_posterior as a .Rdata file.

Value

A list with the following elements:

signature The extracted signature profiles as a matrix; rows are mutation types, columns are samples (e.g. tumors).

exposure The inferred exposures as a matrix of mutation counts; rows are signatures, columns are samples (e.g. tumors).

exposure.p exposure converted to proportions.

multi.chains A `hdpSampleMulti-class` object. This object has the method `chains` which returns a list of `hdpSampleChain-class` objects. Each of these sample chains objects has a method `final_hdpState` (actually the methods seems to be just `hdp`) that returns the `hdpState` from which it was generated.

MultipleSetupAndPosterior

Activate hierarchical Dirichlet processes and run posterior sampling in parallel.

Description

Activate hierarchical Dirichlet processes and run posterior sampling in parallel.

Usage

```
MultipleSetupAndPosterior(
  input.catalog,
  seedNumber = 1,
  K.guess,
  multi.types = FALSE,
  verbose = TRUE,
  post.burnin = 4000,
  post.n = 50,
  post.space = 50,
  post.cpiter = 3,
  post.verbosity = 0,
  CPU.cores = 1,
  num.child.process = 4
)
```

Arguments

<code>input.catalog</code>	Input spectra catalog as a matrix or in ICAMS format.
<code>seedNumber</code>	An integer that is used to generate separate random seeds for each call to dp_activate , and each call of hdp_posterior ; please see the code on how this is done. But repeated calls with same value of <code>seedNumber</code> and other inputs should produce the same results.
<code>K.guess</code>	Initial guess of the number of "raw clusters", which may be larger than the number of signatures (sometimes called "components" in the <code>hdpX</code> code); passed to dp_activate as <code>initcc</code> .
<code>multi.types</code>	<p>A logical scalar or a character vector. If <code>FALSE</code>, The HDP analysis will regard all input spectra as one tumor type.</p> <p>If <code>TRUE</code>, the HDP analysis will infer tumor types based on the string before "::" in their names. e.g. tumor type for "SA.Syn.Ovary-AdenoCA::S.500" would be "SA.Syn.Ovary-AdenoCA"</p> <p>If <code>multi.types</code> is a character vector, then it should be of the same length as the number of columns in <code>input.catalog</code>, and each value is the name of the tumor type of the corresponding column in <code>input.catalog</code>, e.g. <code>c("SA.Syn.Ovary-AdenoCA", "SA.Syn.O</code></p>

verbose	If TRUE then message progress information.
post.burnin	Pass to hdp_posterior burnin.
post.n	Pass to hdp_posterior n.
post.space	Pass to hdp_posterior space.
post.cpiter	Pass to hdp_posterior cpiter.
post.verbosity	Pass to hdp_posterior verbosity.
CPU.cores	Number of CPUs to use; there is no point in making this larger than num.child.process.
num.child.process	Number of posterior sampling chains; can set to 1 for testing.

Value

Invisibly, the clean chlist (output of the [hdp_posterior](#) calls). This is a list of [hdpSampleChain-class](#) objects.

PlotExposure	<i>Plot a single exposure plot</i>
--------------	------------------------------------

Description

Plot a single exposure plot

Usage

```
PlotExposure(exposures, plot.proportion = FALSE, plot.legend = TRUE, ...)
```

Arguments

exposures	Exposures as a numerical matrix (or data.frame) with signatures in rows and samples in columns. Rownames are taken as the signature names and column names are taken as the sample IDs. If you want exp sorted from largest to smallest use SortExp . Do not use column names that start with multiple underscores. The exposures will often be mutation counts, but could also be e.g. mutations per megabase.
plot.proportion	Plot exposure proportions rather than counts.
plot.legend	If TRUE plot a legend.
...	Parameters passed to barplot .

PlotExposureByRange	<i>Plot exposures in multiple plots each with a manageable number of samples.</i>
---------------------	---

Description

Plot exposures in multiple plots each with a manageable number of samples.

Usage

```
PlotExposureByRange(exposures, num.per.line = 30, plot.proportion = FALSE, ...)
```

Arguments

exposures	Exposures as a numerical matrix (or data.frame) with signatures in rows and samples in columns. Rownames are taken as the signature names and column names are taken as the sample IDs. If you want exposures sorted from largest to smallest use SortExp . Do not use column names that start with multiple underscores. The exposures will often be mutation counts, but could also be e.g. mutations per megabase.
num.per.line	Number of samples to show in each plot.
plot.proportion	Plot exposure proportions rather than counts.
...	Other arguments passed to PlotExposure . If ylab is not included, it defaults to a value depending on plot.proportion. If col is not supplied the function tries to do something reasonable.

PrepInit	<i>Initialize hdp object Allocate process index for hdp initialization. Prepare for hdp_init</i>
----------	--

Description

Initialize hdp object Allocate process index for hdp initialization. Prepare for hdp_init

Usage

```
PrepInit(multi.types, input.catalog, verbose, K.guess)
```

Arguments

multi.types	TODO
input.catalog	TODO
verbose	TODO
K.guess	TODO

RunHdpParallel	<i>Extract mutational signatures and optionally compare them to existing signatures and exposures.</i>
----------------	--

Description

Extract mutational signatures and optionally compare them to existing signatures and exposures.

Usage

```
RunHdpParallel(
  input.catalog,
  seedNumber = 1,
  K.guess,
  multi.types = FALSE,
  verbose = TRUE,
  post.burnin = 4000,
  post.n = 50,
  post.space = 50,
  post.cpiter = 3,
  post.verbosity = 0,
  CPU.cores = 1,
  num.child.process = 4,
  cos.merge = 0.9,
  min.sample = 1,
  ground.truth.sig = NULL,
  ground.truth.exp = NULL,
  overwrite = TRUE,
  out.dir = NULL
)
```

Arguments

<code>input.catalog</code>	Input spectra catalog as a matrix or in ICAMS format.
<code>seedNumber</code>	An integer that is used to generate separate random seeds for each call to dp_activate , and each call of hdp_posterior ; please see the code on how this is done. But repeated calls with same value of <code>seedNumber</code> and other inputs should produce the same results.
<code>K.guess</code>	Initial guess of the number of "raw clusters", which may be larger than the number of signatures (sometimes called "components" in the <code>hdp</code> code); passed to dp_activate as <code>initcc</code> .
<code>multi.types</code>	<p>A logical scalar or a character vector. If <code>FALSE</code>, The HDP analysis will regard all input spectra as one tumor type.</p> <p>If <code>TRUE</code>, the HDP analysis will infer tumor types based on the string before ":" in their names. e.g. tumor type for "SA.Syn.Ovary-AdenoCA::S.500" would be "SA.Syn.Ovary-AdenoCA"</p> <p>If <code>multi.types</code> is a character vector, then it should be of the same length as the number of columns in <code>input.catalog</code>, and each value is the name of the tumor type of the corresponding column in <code>input.catalog</code>, e.g. <code>c("SA.Syn.Ovary-AdenoCA", "SA.Syn.O</code></p>

<code>verbose</code>	If TRUE then message progress information.
<code>post.burnin</code>	Pass to hdp_posterior burnin.
<code>post.n</code>	Pass to hdp_posterior n.
<code>post.space</code>	Pass to hdp_posterior space.
<code>post.cpiter</code>	Pass to hdp_posterior cpiter.
<code>post.verbosity</code>	Pass to hdp_posterior verbosity.
<code>CPU.cores</code>	Number of CPUs to use; there is no point in making this larger than <code>num.child.process</code> .
<code>num.child.process</code>	Number of posterior sampling chains; can set to 1 for testing.
<code>cos.merge</code>	The cosine similarity threshold for merging raw clusters from the posterior sampling chains into "components" i.e. signatures; passed to hdp_extract_components .
<code>min.sample</code>	A "component" (i.e. signature) must have at least this many samples; passed to hdp_extract_components .
<code>ground.truth.sig</code>	Optional. Either a string with the path to file with ground truth signatures or and ICAMS catalog with the ground truth signatures. These are the signatures used to construct the ground truth spectra.
<code>ground.truth.exp</code>	Optional. Ground truth exposure matrix or path to file with ground truth exposures. If NULL skip checks that need this information.
<code>overwrite</code>	If TRUE, overwrite <code>out.dir</code> if it is non-NULL and exists.
<code>out.dir</code>	If not NULL, output including data and plots will be saved in <code>out.dir</code> .

Value

Invisibly, a list with the following elements:

signature The extracted signature profiles as a matrix; rows are mutation types, columns are samples (e.g. tumors).

exposure The inferred exposures as a matrix of mutation counts; rows are signatures, columns are samples (e.g. tumors).

multi.chains A [hdpSampleMulti-class](#) object. This object has the method [chains](#) which returns a list of [hdpSampleChain-class](#) objects. Each of these sample chains objects has a method [final_hdpState](#) (actually the methods seems to be just [hdp](#)) that returns the [hdpState](#) from which it was generated.

SetupAndPosterior

Generate an HDP Gibbs sampling chain from a spectra catalog.

Description

Generate an HDP Gibbs sampling chain from a spectra catalog.

Usage

```
SetupAndPosterior(
  input.catalog,
  seedNumber = 1,
  K.guess,
  multi.types = FALSE,
  verbose = TRUE,
  post.burnin = 4000,
  post.n = 50,
  post.space = 50,
  post.cpiter = 3,
  post.verbosity = 0
)
```

Arguments

<code>input.catalog</code>	Input spectra catalog as a matrix or in ICAMS format.
<code>seedNumber</code>	An integer that is used to generate separate random seeds for each call to dp_activate , and each call of hdp_posterior ; please see the code on how this is done. But repeated calls with same value of <code>seedNumber</code> and other inputs should produce the same results.
<code>K.guess</code>	Initial guess of the number of "raw clusters", which may be larger than the number of signatures (sometimes called "components" in the <code>hdpX</code> code); passed to dp_activate as <code>initcc</code> .
<code>multi.types</code>	A logical scalar or a character vector. If <code>FALSE</code> , The HDP analysis will regard all input spectra as one tumor type. If <code>TRUE</code> , the HDP analysis will infer tumor types based on the string before "::" in their names. e.g. tumor type for "SA.Syn.Ovary-AdenoCA::S.500" would be "SA.Syn.Ovary-AdenoCA" If <code>multi.types</code> is a character vector, then it should be of the same length as the number of columns in <code>input.catalog</code> , and each value is the name of the tumor type of the corresponding column in <code>input.catalog</code> , e.g. <code>c("SA.Syn.Ovary-AdenoCA", "SA.Syn.Ovary-AdenoCA")</code>
<code>verbose</code>	If <code>TRUE</code> then message progress information.
<code>post.burnin</code>	Pass to hdp_posterior <code>burnin</code> .
<code>post.n</code>	Pass to hdp_posterior <code>n</code> .
<code>post.space</code>	Pass to hdp_posterior <code>space</code> .
<code>post.cpiter</code>	Pass to hdp_posterior <code>cpiter</code> .
<code>post.verbosity</code>	Pass to hdp_posterior <code>verbosity</code> .

Value

Invisibly, an [hdpSampleChain-class](#) object as returned from [hdp_posterior](#).

SortExp	<i>Sort columns of an exposure matrix from largest to smaller (or vice versa).</i>
---------	--

Description

Sort columns of an exposure matrix from largest to smaller (or vice versa).

Usage

```
SortExp(exposures, decreasing = TRUE)
```

Arguments

exposures	The exposures to sort; columns are samples.
decreasing	If TRUE sort from largest to smallest.

Index

AnalyzeAndPlotretval, [2](#)

barplot, [7](#)

chains, [3](#), [4](#), [6](#), [10](#)

ChainsDiagnosticPlot, [2](#)

CombinePosteriorChains, [3](#)

dp_activate, [5](#), [6](#), [9](#), [11](#)

ExtendIterationAndPosterior, [4](#)

final_hdpState, [3](#), [4](#), [6](#), [10](#)

hdp_extract_components, [4](#), [5](#), [10](#)

hdp_posterior, [5–7](#), [9–11](#)

ICAMS, [2](#), [3](#), [5](#), [6](#), [9–11](#)

MultipleSetupAndPosterior, [6](#)

PlotExposure, [7](#), [8](#)

PlotExposureByRange, [8](#)

PrepInit, [8](#)

RunHdpParallel, [9](#)

SetupAndPosterior, [10](#)

SortExp, [7](#), [8](#), [12](#)