

Package ‘sitePath’

April 15, 2020

Type Package

Title Detection of sites with fixation of amino acid substitutions in protein evolution

Version 1.2.3

Author Chengyang Ji, Aiping Wu

Maintainer Chengyang Ji <chengyang.ji12@alumni.xjtlu.edu.cn>

Description The package does hierarchical search for fixation events given multiple sequence alignment and phylogenetic tree. These fixation events can be specific to a phylogenetic lineages or shared by multiple lineages.

License MIT + file LICENSE

Depends R (>= 3.6.0)

Imports ape, seqinr, Rcpp, methods, graphics, utils, stats

Suggests testthat, knitr, rmarkdown, BiocStyle

LinkingTo Rcpp

RoxygenNote 7.1.0

Encoding UTF-8

VignetteBuilder knitr

URL <https://wuaipinglab.github.io/sitePath/>

BugReports <https://github.com/wuaipinglab/sitePath/issues>

biocViews Alignment, MultipleSequenceAlignment, Software

git_url <https://git.bioconductor.org/packages/sitePath>

git_branch RELEASE_3_10

git_last_commit 86e44d4

git_last_commit_date 2020-04-01

Date/Publication 2020-04-14

R topics documented:

addMSA	2
extractTips	3
findSites	4
h3n2_align	5

h3n2_align_reduced	6
h3n2_tree	6
h3n2_tree_reduced	6
plot.lineagePath	7
plot.sitePath	7
plotSingleSite	8
pre-assessment	9
setSiteNumbering	10
treemer	11
zikh_align	12
zikh_align_reduced	12
zikh_tree	13
zikh_tree_reduced	13

Index 14

addMSA	<i>Prepare data for sitePath analysis</i>
--------	---

Description

sitePath requires both tree and sequence alignment to do the analysis. addMSA wraps read.alignment function in seqinr package and helps match names in tree and sequence alignment. Either provide the file path to an alignment file and its format or an alignment object from the return of read.alignment function. If both the file path and alignment object are given, the function will use the sequence in the alignment file.

Usage

```
addMSA(tree, msaPath = "", msaFormat = "", alignment = NULL)
```

Arguments

tree	a phylo object. This commonly can be from tree parsing function in ape or ggtree. All the tip.label should be found in the sequence alignment.
msaPath	The file path to the multiple sequence alignment file
msaFormat	The format of the multiple sequence alignment file
alignment	an alignment object. This commonly can be from sequence parsing function in the seqinr package. Sequence names in the alignment should include all tip.label in the tree

Value

addMSA returns a phylo object with matched multiple sequence alignment

Examples

```
data(zikh_tree)
msaPath <- system.file('extdata', 'ZIKV.fasta', package = 'sitePath')
addMSA(zikh_tree, msaPath = msaPath, msaFormat = 'fasta')
```

extractTips	<i>Extract info for fixation of a single site</i>
-------------	---

Description

The result of `fixationSites` contains all the possible sites with fixation mutation. The function `extractTips` retrieves the name of the tips involved in the fixation.

The function `extractSite` can be used to extract the fixation info of a single site.

Usage

```
## S3 method for class 'fixationSites'
extractTips(x, site, select = 1, ...)

## S3 method for class 'multiFixationSites'
extractTips(x, site, select = 1, ...)

## S3 method for class 'sitePath'
extractTips(x, select = 1, ...)

## S3 method for class 'fixationSites'
extractSite(x, site, ...)

## S3 method for class 'multiFixationSites'
extractSite(x, site, ...)
```

Arguments

<code>x</code>	A <code>fixationSites</code> or a <code>multiFixationSites</code> or a <code>sitePath</code> object.
<code>site</code>	A site predicted to experience fixation.
<code>select</code>	For a site, there theoretically might be more than one fixation on different lineages. You may use this argument to extract for a specific fixation of a site. The default is the first fixation of the site.
<code>...</code>	Other arguments

Value

The function `extractTips` returns the name of the tips involved in the fixation.

The function `extractSite` returns a `sitePath` object

Examples

```
data(zikv_tree_reduced)
data(zikv_align_reduced)
tree <- addMSA(zikv_tree_reduced, alignment = zikv_align_reduced)
mutations <- fixationSites(lineagePath(tree))
extractTips(mutations, 139)
extractSite(mutations, 139)
```

Description

Single nucleotide polymorphism (SNP) in the whole package refers to variation of amino acid. `findSNPsite` will try to find SNP in the multiple sequence alignment. A reference sequence and gap character may be specified to number the site. This is irrelevant to the intended analysis but might be helpful to evaluate the performance of `fixationSites`.

After finding the `lineagePath` of a phylogenetic tree, `fixationSites` uses the result to find those sites that show fixation on some, if not all, of the lineages. Parallel evolution is relatively common in RNA virus. There is chance that some site be fixed in one lineage but does not show fixation because of different sequence context.

After finding the `lineagePath` of a phylogenetic tree, `multiFixationSites` uses random sampling on the original tree and applies the method used in `fixationSites` to each sampled tree and summarize the results from all the samples.

Usage

```
SNPsites(tree, minSNP = NULL)

## S3 method for class 'lineagePath'
fixationSites(
  paths,
  minEffectiveSize = NULL,
  searchDepth = 1,
  method = c("compare", "insert", "delete"),
  ...
)

## S3 method for class 'lineagePath'
multiFixationSites(
  paths,
  samplingSize = NULL,
  samplingTimes = 100,
  minEffectiveSize = 0,
  searchDepth = 1,
  method = c("compare", "insert", "delete"),
  ...
)
```

Arguments

<code>tree</code>	The return from <code>addMSA</code> function
<code>minSNP</code>	Minimum number of amino acid variation to be a SNP
<code>paths</code>	a <code>lineagePath</code> object returned from <code>lineagePath</code> function
<code>minEffectiveSize</code>	A vector of two integers to specify minimum tree tips involved before/after mutation. Otherwise the mutation will not be counted into the return. If more

	than one number is given, the ancestral takes the first and descendant takes the second as the minimum. If only given one number, it's the minimum for both ancestral and descendant.
searchDepth	The function uses heuristic search but the termination of the search cannot be intrinsically decided. searchDepth is needed to tell the search when to stop.
method	The strategy for predicting the fixation. The basic approach is entropy minimization and can be achieved by adding or removing fixation point, or by comparing the two.
...	further arguments passed to or from other methods.
samplingSize	The number of tips sampled for each round of resampling. It should be at least 10th and at most nine 10ths of the tree tips.
samplingTimes	The total times of random sampling to do. It should be greater than 100.

Value

SNPsite returns a list of qualified SNP site

fixationSites returns a list of fixation mutations with names of the tips involved.

multiFixationSites returns sites with multiple fixations.

Examples

```
data(zikv_tree_reduced)
data(zikv_align_reduced)
tree <- addMSA(zikv_tree_reduced, alignment = zikv_align_reduced)
SNPsites(tree)
fixationSites(lineagePath(tree))
```

h3n2_align

Multiple sequence alignment of H3N2's HA protein

Description

The raw protein sequences were downloaded from NCBI database.

Usage

```
data(h3n2_align)
```

Format

a alignment object

h3n2_align_reduced *Truncated data for runnable example*

Description

This is a truncated version of [h3n2_align](#)

Usage

```
data(h3n2_align_reduced)
```

Format

a alignment object

h3n2_tree *Phylogenetic tree of H3N2's HA protein*

Description

Tree was built from [h3n2_align](#) using RAxML with default settings.

Usage

```
data(h3n2_tree)
```

Format

a phylo object

h3n2_tree_reduced *Truncated data for runnable example*

Description

This is a truncated version of [h3n2_tree](#)

Usage

```
data(h3n2_tree_reduced)
```

Format

a phylo object

plot.lineagePath	<i>Visualize phylogenetic lineages</i>
------------------	--

Description

Visualize `lineagePath` object. A tree diagram will be plotted and paths are black solid line while the trimmed nodes and tips will use grey dashed line.

Usage

```
## S3 method for class 'lineagePath'
plot(x, y = TRUE, showTips = FALSE, ...)
```

Arguments

x	A <code>lineagePath</code> object
y	Whether plot the nodes from the extendedSearch in <code>fixationSites</code>
showTips	Whether to plot the tip labels. The default is FALSE.
...	Arguments in <code>plot.phylo</code> functions.

Value

The function only makes plot and returns no value (It behaviors like the generic `plot` function).

Examples

```
data(zikv_tree)
data(zikv_align)
tree <- addMSA(zikv_tree, alignment = zikv_align)
plot(lineagePath(tree))
```

plot.sitePath	<i>Plot the fixation mutation</i>
---------------	-----------------------------------

Description

Visualize the `sitePath` object which is the basic unit of the result of `fixationSites` and `multiFixationSites`.

Usage

```
## S3 method for class 'sitePath'
plot(x, y = NULL, showTips = FALSE, ...)
```

Arguments

x	A <code>sitePath</code> object
y	A <code>sitePath</code> object can have more than one fixation path. This is to select which path to plot. The default is NULL which will plot all the paths.
showTips	Whether to plot the tip labels. The default is FALSE.
...	Arguments in <code>plot.phylo</code> functions and other arguments.

Value

The function only makes plot and returns no value (It behaviors like the generic `plot` function).

See Also

[plotSingleSite](#)

Examples

```
data(zikv_align_reduced)
data(zikv_tree_reduced)
tree <- addMSA(zikv_tree_reduced, alignment = zikv_align_reduced)
paths <- lineagePath(tree)
fixations <- fixationSites(paths)
plot(fixations[[1]])
```

<code>plotSingleSite</code>	<i>Color the tree by a single site</i>
-----------------------------	--

Description

For `lineagePath`, the tree will be colored according to the amino acid of the site. The color scheme tries to assign distinguishable color for each amino acid.

For `fixationSites`, it will color the ancestral tips in red, descendant tips in blue and excluded tips in grey.

For `multiFixationSites`, it will color the tips which have their site fixed. The color will use the same amino acid color scheme as `plotSingleSite.lineagePath`

Usage

```
## S3 method for class 'lineagePath'
plotSingleSite(x, site, showPath = FALSE, showTips = FALSE, ...)

## S3 method for class 'fixationSites'
plotSingleSite(x, site, select = NULL, ...)

## S3 method for class 'multiFixationSites'
plotSingleSite(x, site, select = NULL, ...)
```

Arguments

<code>x</code>	A <code>fixationSites</code> object from fixationSites or the return from addMSA function.
<code>site</code>	One of the mutations in the <code>fixationSites</code> object. It should be from the <code>names</code> of the object. Or an integer to indicate a site could be provide. The numbering is consistent with the reference defined at fixationSites .
<code>showPath</code>	If plot the lineage result from <code>lineagePath</code> .
<code>showTips</code>	Whether to plot the tip labels. The default is <code>FALSE</code> .
<code>...</code>	Arguments in <code>plot.phylo</code> functions and other arguments.
<code>select</code>	Select which fixation path in to plot. The default is <code>NULL</code> which will plot all the fixations.

Value

The function only makes plot and returns no value (It behaviors like the generic `plot` function).

See Also

[plot.sitePath](#)

Examples

```
data(zikv_tree)
data(zikv_align)
tree <- addMSA(zikv_tree, alignment = zikv_align)
paths <- lineagePath(tree)
plotSingleSite(paths, 139)
fixations <- fixationSites(paths)
plotSingleSite(fixations, 139)
## Not run:
multiFixations <- multiFixationSites(paths)
plotSingleSite(multiFixations, 1542)

## End(Not run)
```

```
pre-assessment
```

Things can be done before the analysis

Description

`similarityMatrix` calculates similarity between aligned sequences The similarity matrix can be used in [groupTips](#) or [lineagePath](#)

`sneakPeek` is intended to plot 'similarity' (actually the least percentage of 'major SNP') as a threshold against number of output `lineagePath`. This plot is intended to give user a rough view about how many lineages they could expect from the 'similarity' threshold in the function [lineagePath](#). The number of `lineagePath` is preferably not be too many or too few. The result excludes where the number of `lineagePath` is greater than number of tips divided by 20 or user-defined `maxPath`. The zero `lineagePath` result will also be excluded.

Usage

```
similarityMatrix(tree)

sneakPeek(tree, step = 10, maxPath = NULL, minPath = 1, makePlot = FALSE)
```

Arguments

<code>tree</code>	The return from addMSA function
<code>step</code>	the 'similarity' window for calculating and plotting. To better see the impact of threshold on path number. The default is 10.
<code>maxPath</code>	maximum number of path to return show in the plot. The number of path in the raw tree can be far greater than trimmed tree. To better see the impact of threshold on path number. This is preferably specified. The default is one 20th of tree tip number.

minPath	minimum number of path to return show in the plot. To better see the impact of threshold on path number. The default is 1.
makePlot	whether make a dot plot when return

Value

similarityMatrix returns a diagonal matrix of similarity between sequences

sneakPeek return the similarity threshold against number of lineagePath. There will be a simple dot plot between threshold and path number if makePlot is TRUE.

Examples

```
data('zikkv_tree')
data('zikkv_align')
tree <- addMSA(zikkv_tree, alignment = zikkv_align)
simMatrix <- similarityMatrix(tree)
sneakPeek(tree)
```

setSiteNumbering *Set site numbering to the reference sequence*

Description

A reference sequence can be used to define a global site numbering scheme for multiple sequence alignment. The gap in the reference will be skipped so the site ignored in numbering.

Usage

```
## S3 method for class 'phylo'
setSiteNumbering(x, reference = NULL, gapChar = "-", ...)

## S3 method for class 'lineagePath'
setSiteNumbering(x, reference = NULL, gapChar = "-", ...)
```

Arguments

x	The object to set site numbering. It could be a phylo object after <code>addMSA</code> or a lineagePath object. The function for <code>fixaitonSites</code> and <code>multiFixationSites</code> will be added in later version.
reference	Name of reference for site numbering. The name has to be one of the sequences' name. The default uses the intrinsic alignment numbering
gapChar	The character to indicate gap. The numbering will skip the gapChar for the reference sequence.
...	further arguments passed to or from other methods.

Value

A phylo object with site numbering mapped to reference sequence

Examples

```
data(zikv_tree)
msaPath <- system.file('extdata', 'ZIKV.fasta', package = 'sitePath')
tree <- addMSA(zikv_tree, msaPath = msaPath, msaFormat = 'fasta')
setSiteNumbering(tree)
```

treemer

Topology-dependent tree trimming

Description

groupTips uses sequence similarity to group tree tips. Members in a group are always constrained to share the same ancestral node. Similarity between two tips is derived from their multiple sequence alignment. The site will not be counted into total length if both are gap. Similarity is calculated as number of matched divided by the corrected total length. So far the detection of divergence is based on one simple rule: the minimal pairwise similarity. The two branches are decided to be divergent if the similarity is lower than the threshold. (Other more statistical approaches such as Kolmogorov-Smirnov Tests among pair-wise distance could be introduced in the future)

lineagePath finds the lineages of a phylogenetic tree providing the corresponding sequence alignment. This is done by finding 'major SNPs' which usually accumulate along the evolutionary pathways. are added.

Usage

```
groupTips(
  tree,
  similarity = NULL,
  simMatrix = NULL,
  forbidTrivial = TRUE,
  tipnames = TRUE
)
```

```
lineagePath(tree, similarity = NULL, simMatrix = NULL, forbidTrivial = TRUE)
```

Arguments

tree	The return from addMSA function
similarity	Similarity threshold for tree trimming in groupTips. If not provided, the mean similarity subtract standard deviation of all sequences will be used. And for lineagePath, this decides how minor SNPs are to remove. If provided as fraction between 0 and 1, then the minimum number of SNP will be total tips times similarity. If provided as integer greater than 1, the minimum number will be similarity. The default similarity is 0.1 for lineagePath.
simMatrix	A diagonal matrix of similarities for each pair of sequences. This parameter will not have effect in the function lineagePath.
forbidTrivial	Does not allow trivial trimming
tipnames	If return as tipnames

Value

grouping of tips
path represent by node number

Examples

```
data('zikh_tree')  
data('zikh_align')  
tree <- addMSA(zikh_tree, alignment = zikh_align)  
groupTips(tree, 0.996)  
lineagePath(tree)
```

zikh_align	<i>Multiple sequence alignment of Zika virus polyprotein</i>
------------	--

Description

The raw protein sequences were downloaded from ViPR database (<https://www.viprbrc.org/>) and aligned using MAFFT. with default settings.

Usage

```
data(zikh_align)
```

Format

a alignment object

zikh_align_reduced	<i>Truncated data for runnable example</i>
--------------------	--

Description

This is a truncated version of [zikh_align](#)

Usage

```
data(zikh_align_reduced)
```

Format

a alignment object

`zikh_tree`*Phylogenetic tree of Zika virus polyprotein*

Description

Tree was built from [zikh_align](#) using RAxML with default settings. The tip ANK57896 was used as outgroup to root the tree.

Usage

```
data(zikh_tree)
```

Format

a phylo object

`zikh_tree_reduced`*Truncated data for runnable example*

Description

This is a truncated version of [zikh_tree](#)

Usage

```
data(zikh_tree_reduced)
```

Format

a phylo object

Index

*Topic **datasets**

- h3n2_align, 5
- h3n2_align_reduced, 6
- h3n2_tree, 6
- h3n2_tree_reduced, 6
- zikv_align, 12
- zikv_align_reduced, 12
- zikv_tree, 13
- zikv_tree_reduced, 13

addMSA, 2, 4, 8–11

extractSite (extractTips), 3
extractTips, 3

findSites, 4
fixationSites, 3, 7, 8
fixationSites (findSites), 4

groupTips, 9
groupTips (treemer), 11

h3n2_align, 5, 6
h3n2_align_reduced, 6
h3n2_tree, 6, 6
h3n2_tree_reduced, 6

lineagePath, 4, 7, 9
lineagePath (treemer), 11

multiFixationSites, 7
multiFixationSites (findSites), 4

names, 8

plot, 7–9
plot.lineagePath, 7
plot.sitePath, 7, 9
plotSingleSite, 8, 8
pre-assessment, 9

setSiteNumbering, 10
similarityMatrix (pre-assessment), 9
sneakPeek (pre-assessment), 9
SNPsites (findSites), 4

treemer, 11

zikv_align, 12, 12, 13
zikv_align_reduced, 12
zikv_tree, 13, 13
zikv_tree_reduced, 13