Package 'longreadvqs'

October 20, 2025

Title Viral Quasispecies Comparison from Long-Read Sequencing Data

Version 0.1.4

Description Performs variety of viral quasispecies diversity analyses [see Pamorn-chainavakul et al. (2024) <doi:10.21203/rs.3.rs-4637890/v1>] based on long-read sequence alignment. Main functions include 1) sequencing error and other noise minimization and read sampling, 2) Single nucleotide variant (SNV) profiles comparison, and 3) viral quasispecies profiles comparison and visualization.

License GPL-3

URL https://github.com/NakarinP/longreadvqs

BugReports https://github.com/NakarinP/longreadvqs/issues

Encoding UTF-8

RoxygenNote 7.3.3

Imports ape, Biostrings, cowplot, dplyr, ggplot2, ggpubr, grDevices, IRanges, magrittr, methods, plyr, purrr, pwalign, QSutils, RColorBrewer, reshape2, scales, seqinr, stats, stringdist, stringr, tibble, tidyr

Depends R (>= 4.4)

Suggests knitr, rmarkdown, testthat (>= 3.0.0)

VignetteBuilder knitr

Config/testthat/edition 3

NeedsCompilation no

Author Nakarin Pamornchainavakul [aut, cre] (ORCID: https://orcid.org/0000-0003-0378-0316)

Repository CRAN

Date/Publication 2025-10-20 18:20:07 UTC

2 AAcompare

Contents

AAcompare	2
iltfast	3
gapremove	4
otucompare	5
octopt	
snvcompare	
/qsassess	7
yqscompare	9
graphy of graphs of the state o	(
yqsout	2
yqsresub	
7qssub	4
1	

AAcompare

Index

Comparing viral quasispecies diversity metrics at amino acid level

Description

Pools noise-minimized down-sampled read samples and compares their diversity metrices based on protein haplotype and single amino acid variation (SAV) group that is classified by k-means clustering of SAV distance. This function is a subset of "vqscompare" function.

Usage

```
AAcompare(
  samplelist = list(BC1, BC2, BC3),
  kmeans.n = 20,
  removestopcodon = FALSE
)
```

Arguments

samplelist List of samples, i.e., name of resulting objects from "vqsassess" or "vqscustom-

pct" functions, for example list(BC1, BC2, BC3).

kmeans.n Number of single amino acid variation (SAV) groups needed from k-means clus-

tering on multidimensional scale (MDS) of all samples' pairwise SAV distance.

removestopcodon

Remove the last amino acid (expected to be a stop codon) from translated amino acid sequences before further analysis (optional). If not specified or if removestop-codon = FALSE, the last amino acid will not be removed (default).

filtfast 3

Value

List of 1) "aadiv": comparative table of viral quasispecies diversity metrics between listed samples based on translated reads calculated by QSutils package, and 2) "savgrpdiv": comparative table of single amino acid (SAV) group diversity metrics between listed samples calculated from consensus amino acid sequence of each SAV group

Examples

```
## Locate input FASTA files------
sample1filepath <- system.file("extdata", "s1.fasta", package = "longreadvqs")
sample2filepath <- system.file("extdata", "s2.fasta", package = "longreadvqs")

## Prepare data for viral quasispecies comparison between two samples------
set.seed(123)
sample1 <- vqsassess(sample1filepath, pct = 5, samsize = 50, label = "sample1")
sample2 <- vqsassess(sample2filepath, pct = 5, samsize = 50, label = "sample2")

## Compare protein haplotype and SAV group (4 clusters) diversity metrics between two samples-----
AAcompare(samplelist = list(sample1, sample2), kmeans.n = 4)</pre>
```

filtfast

Filtering highly dissimilar reads/sequences out of the alignment

Description

Removes reads/sequences of which Hamming similarity to the consensus of all reads/sequences in the alignment is less than the specified quantile (qt) of the similarity distribution.

Usage

```
filtfast(fasta, qt = 0.25, fastaname = "filteredfast.fasta")
```

Arguments

fasta Input as a read or multiple sequence alignment in FASTA format

qt If Hamming similarity score of a read/sequence to the consensus of all reads/sequences is less than the specified quantile (qt) of the similarity distribution, that read/sequence will be removed.

fastaname Output file name in FASTA format

Value

FASTA read or multiple sequence alignment written out to the input directory

4 gapremove

Examples

```
## Locate input FASTA file------
fastafilepath <- system.file("extdata", "dissimfast.fasta", package = "longreadvqs")

## Indicate output directory and file name------
outfast <- tempfile()

## Remove reads/sequences that the similarity < 1st quartile (0.25 quantile)------
filtfast(fastafilepath, qt = 0.25, fastaname = outfast)</pre>
```

gapremove

Removing gap-rich positions and/or reads/sequences

Description

Removes nucleotide positions (vertical) and/or reads/sequences (horizontal) that contain gaps more than the specified cut-off percentage from the alignment.

Usage

```
gapremove(fasta, vgappct = 70, hgappct = 70, fastaname = "filteredfast.fasta")
```

Arguments

fasta Input as a read or multiple sequence alignment in FASTA format

vgappct The percent cut-off of vertical gap (-), i.e., if a position in the alignment has

%gap >= vgappct, that position will be removed.

hgappet The percent cut-off of horizontal gap (-), i.e., if a sequence or read in the align-

ment has %gap >= hgappet, that sequence or read will be removed.

fastaname Output file name in FASTA format

Value

FASTA read or multiple sequence alignment written out to the input directory

```
## Locate input FASTA file------
fastafilepath <- system.file("extdata", "gaprichfast.fasta", package = "longreadvqs")

## Indicate output directory and file name------
outfast <- tempfile()

## Remove positions with gap >= 60% and reads/sequences with gap >= 10%------
gapremove(fastafilepath, vgappct = 60, hgappct = 10, fastaname = outfast)
```

otucompare 5

otucompare	Comparing operational taxonomic unit (OTU) by k-means clustering between samples

Description

Pools noise-minimized down-sampled read samples and compares their diversity based on operational taxonomic unit (OTU) classified by k-means clustering of single nucleotide variant (SNV) distance. This function is a subset of "vqscompare" function.

Usage

```
otucompare(samplelist = list(BC1, BC2, BC3), kmeans.n = 20)
```

Arguments

samplelist List of samples, i.e., name of resulting objects from "vqsassess" or "vqscustom-

pct" functions, for example list(BC1, BC2, BC3).

Number of operational taxonomic units (OTUs) needed from k-means clustering

on multidimensional scale (MDS) of all samples' pairwise genetic distance.

Value

Comparative table of OTU diversity metrics between listed samples calculated from consensus sequence of each OTU by QSutils package

6 pctopt

pctopt Optimizing cut-off percentage for noise minimization

Description

Finds an optimal cut-off percentage for noise minimization (in vqssub, vqsassess, and vqscustompct functions) that can decrease the number of singleton haplotypes to less than the desired percentage of the total reads.

Arguments

fasta	Input as a read alignment in FASTA format	
pctsing	The desired percentage of singleton haplotypes relative to the total reads in the alignment.	
method	Sequencing error and noise minimization methods that replace low frequency nucleotide base (less than the "pct" cut-off) with consensus base of that position ("conbase": default) or with base of the dominant haplotype ("domhapbase").	
samplingfirst	Downsampling before (TRUE) or after (FALSE: default) the noise minimization.	
gappct	The percent cut-off particularly specified for gap (-). If it is not specified or less than "pct", "gappct" will be equal to "pct" (default).	
ignoregappositions		
	Replace all nucleotides in the positions in the alignment containing gap(s) with gap. This will make such positions no longer single nucleotide variant (SNV). The default is "FALSE".	
samsize	Sample size (number of reads) after down-sampling. If it is not specified or more than number of reads in the original alignment, down-sampling will not be performed (default).	
label	String within quotation marks indicating name of read alignment (optional).	

Value

An optimal cut-off percentage for noise minimization of an input sample and parameter settings. If label is specified, the output will be a data frame with percentage of singleton haplotypes at each cut-off percentage from zero to the optimal cut-off percentage.

```
## Locate input FASTA file------
fastafilepath <- system.file("extdata", "s1.fasta", package = "longreadvqs")

## Find an cut-off percentage that creates singleton haplotypes less than 50% of the alignment.----
pctopt(fastafilepath, pctsing = 50, label = "s1")</pre>
```

snvcompare 7

snvcompare	Plotting single nucleotide variant (SNV) frequency in read alignment across different samples

Description

Compares single nucleotide variant (SNV) profile between noise-minimized down-sampled read samples using cowplot's "plot_grid" function. The resulting plot may help evaluating the optimal cut-off percentage of low frequency nucleotide base used in "vqsassess", "vqscustompct", or "vqs-sub" functions.

Arguments

samplelist List of samples, i.e., name of resulting objects from "vqsassess" or "vqscustom-

pct" functions, for example list(BC1, BC2, BC3).

ncol Number of columns for multiple plots (see cowplot's "plot_grid" function)

Value

Comparative plot of SNV frequency in read alignment across different samples

Examples

vqsassess Sequencing error and noise minimization, read down-sampling, and data preparation for viral quasispecies comparison

Description

Minimizes potential long-read sequencing error and noise based on the specified cut-off percentage of low frequency nucleotide base and down-samples read for further comparison with other samples. The output of this function is a list of several objects representing diversity of each sample that must be used as an input for other functions such as "snvcompare" or "vqscompare".

8 vqsassess

Arguments

fasta

method Sequencing error and noise minimization methods that replace low frequency

Input as a read alignment in FASTA format

nucleotide base (less than the "pct" cut-off) with consensus base of that position ("conbase": default) or with base of the dominant haplotype ("domhapbase").

samplingfirst Downsampling before (TRUE) or after (FALSE: default) the noise minimiza-

tion.

pct Percent cut-off defining low frequency nucleotide base that will be replaced

(must be specified).

gappct The percent cut-off particularly specified for gap (-). If it is not specified or less

than "pct", "gappct" will be equal to "pct" (default).

ignoregappositions

Replace all nucleotides in the positions in the alignment containing gap(s) with gap. This will make such positions no longer single nucleotide variant (SNV).

The default is "FALSE".

samsize Sample size (number of reads) after down-sampling. If it is not specified or

more than number of reads in the original alignment, down-sampling will not be

performed (default).

label String within quotation marks indicating name of read alignment (optional).

Please don't use underscore (_) in the label.

Value

List of 1) "dat": viral quasispecies diversity metrics calculated by QSutils package (similar to "vqs-sub" function's output), 2) "snvhap": SNV profile of each haplotype with frequency and new label for "vqscompare" function, 3) "snv": plot of SNV frequency for "snvcompare" function, 4) "hapre": DNAStringSet of read alignment of each haplotype for "vqscompare" function, 5) "lab": name of sample or read alignment

vqscompare 9

vqscompare	Comparing viral quasispecies profile and operational taxonomic unit (OTU) classified by k-means clustering between samples
	(010) classified by k-means clastering between samples

Description

Pools noise-minimized down-sampled read samples and compares their diversity by 1) viral quasispecies profile (haplotype and metrics from QSutils package), 2) operational taxonomic unit (OTU) classified by k-means clustering of single nucleotide variant (SNV) distance, and 3) visualization of different comparative method, i.e., haplotype, OTU, phylogenetic tree, MDS plot. Such comparisons can also be performed at the amino acid level (protein haplotype and single amino acid variation (SAV) group).

Arguments

samplelist	List of samples, i.e., name of resulting objects from "vqsassess" or "vqscustompct" functions, for example list(BC1, BC2, BC3).
lab_name	Name of variable or type of sample for instance "barcode", "sample", "dpi", or "isolate" (optional).
kmeans.n	Number of operational taxonomic units (OTUs) and single amino acid variation (SAV) groups needed from k-means clustering on multidimensional scale (MDS) of all samples' pairwise SNV and SAV distances.
showhap.n	Number of largest haplotypes (default = 30) labeled in the top five OTUs' MDS plot (optional).
proteincoding	Translate gene or protein-coding reads into amino acid sequences and regroup them into protein haplotypes and single amino acid variation (SAV) groups, which are comparable to haplotypes and OTUs at the nucleotide level, respectively (optional). If not specified or if proteincoding = FALSE, gene translation and downstream analyses will not be performed (default).

removestopcodon

Remove the last amino acid (expected to be a stop codon) from translated amino acid sequences before further analysis (optional). If not specified or if removestop-codon = FALSE, the last amino acid will not be removed (default).

Value

List of 1) "hapdiv": comparative table of viral quasispecies diversity metrics between listed samples calculated by QSutils package, 2) "otudiv": comparative table of operational taxonomic unit (OTU) diversity metrics between listed samples calculated from consensus sequence of each OTU (similar to "otucompare" function's output), 3) "sumsnv_hap_otu": frequency and SNV profile (by position in the alignment) of all haplotypes and OTUs, 4) "fullseq": complete read sequence of all haplotypes, 5) "fulldata": complete read sequence of all haplotypes in every sample with frequency and OTU classification, 6) "summaryplot": visualization of viral quasispecies comparison between samples including 6.1) "happlot": proportion of haplotypes (top left), 6.2) "otuplot": proportion of OTUs (bottom left), 6.3) multidimensional scale (MDS) plots (right) of k-means OTU

10 vqscustompct

("top5otumds": 5 largest groups with major haplotypes labeled and "allotumds": all groups), 7) "aadiv": comparative table of viral quasispecies diversity metrics between listed samples based on translated reads calculated by QSutils package (similar to one of "AAcompare" function's outputs), 8) "savgrpdiv": comparative table of single amino acid (SAV) group diversity metrics between listed samples calculated from consensus amino acid sequence of each SAV group (similar to one of "AAcompare" function's outputs), 9) "sumsav_phap_savgrp": frequency and SAV profile (by position in the alignment) of all protein haplotypes and SAV groups, 10) "fullseq_aa": complete amino acid sequence of all protein haplotypes, 11) "fulldata_aa": complete amino acid sequence of all protein haplotypes in every sample with frequency and SAV group classification, 12) "summaryplot_aa": visualization of viral quasispecies comparison between samples based on translated reads including 12.1) "phapplot": proportion of protein haplotypes (top left), 12.2) "savgrpplot": proportion of SAV groups (bottom left), 12.3) multidimensional scale (MDS) plots (right) of k-means SAV group ("top5savgrpmds": 5 largest SAV groups with major protein haplotypes labeled and "allsavgrpmds": all SAV groups), 7) to 12) will be generated only when proteincoding = TRUE

Examples

vqscustompct

Sequencing error and noise minimization with customized % cut-off at particular nucleotide region, read down-sampling, and data preparation for viral quasispecies comparison

Description

Minimizes potential long-read sequencing error and noise based on the specified cut-off percentages of low frequency nucleotide base and down-samples read for further comparison with other samples. In this function, the cut-off percentage can be specifically adjusted for different ranges of nucleotide positions which is very useful when sequencing error heavily occurs in a particular part of reads. The output of this function is a list of several objects representing diversity of each sample that must be used as an input for other functions such as "snvcompare" or "vqscompare".

vqscustompct 11

Arguments

fasta Input as a read alignment in FASTA format

method Sequencing error and noise minimization methods that replace low frequency

nucleotide base (less than the "pct" cut-off) with consensus base of that position ("conbase": default) or with base of the dominant haplotype ("domhapbase").

samplingfirst Downsampling before (TRUE) or after (FALSE: default) the noise minimiza-

tion.

pct Percent cut-off defining low frequency nucleotide base that will be replaced

(must be specified).

brkpos Ranges of nucleotide positions with different % cut-off specified in "lspct" for

example c("1:50","51:1112") meaning that the first and the second ranges are

nucleotide positions 1 to 50 and 51 to 1112, respectively.

1spct List of customized % cut-off applied to nucleotide ranges set in "brkpos" for

example c(15,8) meaning that 15% and 8% cut-offs will be applied to the first

and the second ranges, respectively.

gappet The percent cut-off particularly specified for gap (-). If it is not specified or less

than "pct", "gappct" will be equal to "pct" (default).

ignoregappositions

Replace all nucleotides in the positions in the alignment containing gap(s) with

gap. This will make such positions no longer single nucleotide variant (SNV).

The default is "FALSE".

samsize Sample size (number of reads) after down-sampling. If it is not specified or

more than number of reads in the original alignment, down-sampling will not be

performed (default).

label String within quotation marks indicating name of read alignment (optional).

Please don't use underscore (_) in the label.

Value

List of 1) "dat": viral quasispecies diversity metrics calculated by QSutils package (similar to "vqs-sub" function's output), 2) "snvhap": SNV profile of each haplotype with frequency and new label for "vqscompare" function, 3) "snv": plot of SNV frequency for "snvcompare" function, 4) "hapre": DNAStringSet of read alignment of each haplotype for "vqscompare" function, 5) "lab": name of sample or read alignment

12 vqsout

Use "snvcompare" function to check whether SNV profile looks better or not-----snvcompare(samplelist = list(nocustom, custom), ncol = 1)

vqsout

Exporting viral quasispecies profile comparison results

Description

Writes out resulting objects from "vqscompare" function as tables (TSV files) and alignment (FASTA file) to the working directory.

Usage

```
vqsout(vqscompare.obj, directory = "path/to/directory")
```

Arguments

vqscompare.obj A resulting object from "vqscompare" function.

directory

Path to desired directory (location) for output files. If it is not specified, the directory will be the current working directory.

Value

TSV files of viral quasispecies profile comparison results and FASTA file of unique haplotype alignment.

vqsresub 13

vqsresub	Computing viral quasispecies diversity metrics of noise-minimized re- peatedly down-sampled read alignments

Description

Minimizes potential long-read sequencing error and noise based on the specified cut-off percentage of low frequency nucleotide base and repeatedly down-samples read for sensitivity analysis of the diversity metrics varied by different sample sizes. The output of this function is a summary of viral quasispecies diversity metrics per each iteration of down-sampling calculated by QSutils package's functions. This function is an extension of "vqssub" function.

Arguments

fasta	Input as a read alignment in FASTA format	
iter	Number of iterations for downsampling after noise minimization.	
method	Sequencing error and noise minimization methods that replace low frequency nucleotide base (less than the "pct" cut-off) with consensus base of that position ("conbase": default) or with base of the dominant haplotype ("domhapbase").	
pct	Percent cut-off defining low frequency nucleotide base that will be replaced (must be specified).	
gappct	The percent cut-off particularly specified for gap (-). If it is not specified or less than "pct", "gappet" will be equal to "pct" (default).	
ignoregappositions		
	Replace all nucleotides in the positions in the alignment containing gap(s) with gap. This will make such positions no longer single nucleotide variant (SNV). The default is "FALSE".	
samsize	Sample size (number of reads) after down-sampling. If it is not specified or more than number of reads in the original alignment, down-sampling will not be performed (default).	
label	String within quotation marks indicating name of read alignment (optional).	

Value

Data frame containing all viral quasispecies diversity metrics calculated by QSutils package, noise minimization, and down-sampling information per each downsampling iteration.

```
## Locate input FASTA file------
fastafilepath <- system.file("extdata", "s1.fasta", package = "longreadvqs")

## Summarize viral quasispecies diversity metrics from five downsampling iterations.------
vqsresub(fastafilepath, iter = 5, pct = 10, samsize = 20, label = "sample1")</pre>
```

14 vqssub

	vqssub	Computing viral quasispecies diversity metrics of noise-minimized down-sampled read alignment
--	--------	---

Description

Minimizes potential long-read sequencing error and noise based on the specified cut-off percentage of low frequency nucleotide base and down-samples read for further comparison with other samples. The output of this function is a summary of viral quasispecies diversity metrics calculated by QSutils package's functions. This function is a subset of "vqsassess" function.

Arguments

fasta	Input as a read alignment in FASTA format	
method	Sequencing error and noise minimization methods that replace low frequency nucleotide base (less than the "pct" cut-off) with consensus base of that position ("conbase": default) or with base of the dominant haplotype ("domhapbase").	
samplingfirst	Downsampling before (TRUE) or after (FALSE: default) the noise minimization.	
pct	Percent cut-off defining low frequency nucleotide base that will be replaced (must be specified).	
gappct	The percent cut-off particularly specified for gap (-). If it is not specified or less than "pct", "gappet" will be equal to "pct" (default).	
ignoregappositions		
	Replace all nucleotides in the positions in the alignment containing gap(s) with gap. This will make such positions no longer single nucleotide variant (SNV). The default is "FALSE".	
samsize	Sample size (number of reads) after down-sampling. If it is not specified or more than number of reads in the original alignment, down-sampling will not be performed (default).	
label	String within quotation marks indicating name of read alignment (optional).	

Value

Data frame containing all viral quasispecies diversity metrics calculated by QSutils package, noise minimization, and down-sampling information.

```
## Locate input FASTA file------
fastafilepath <- system.file("extdata", "s1.fasta", package = "longreadvqs")

## Summarize viral quasispecies diversity metrics------
# From noise-minimized unsampled reads (10% cut-off):
vqssub(fastafilepath, pct = 10, label = "sample1")
# From noise-minimized sampled reads (n = 20):</pre>
```

vqssub 15

```
vqssub(fastafilepath, pct = 10, samsize = 20, label = "sample1")
# From noise-minimized sampled reads with 50% cut-off for gap:
vqssub(fastafilepath, pct = 10, gappct = 50, samsize = 20, label = "sample1")
# From noise-minimized sampled reads but ignore positions with gap:
vqssub(fastafilepath, pct = 10, ignoregappositions = TRUE, samsize = 20, label = "sample1")
# From reads that were down-sampled before noise minimization:
vqssub(fastafilepath, pct = 10, samplingfirst = TRUE, samsize = 20, label = "sample1")
```

Index

```
AAcompare, 2
filtfast, 3
gapremove, 4
otucompare, 5
pctopt, 6
snvcompare, 7
vqsassess, 7
vqscompare, 9
vqscustompct, 10
vqsout, 12
vqsresub, 13
vqssub, 14
```