

Package ‘MHCtools’

August 11, 2019

Type Package

Title Analysis of MHC Data in Non-Model Species

Version 1.2.1

Description Ten tools for analysis of major histocompatibility complex (MHC) data in non-model species. The functions are tailored for amplicon data sets that have been filtered using the 'dada2' method (for more information visit <https://benjjneb.github.io/dada2>), but even other data sets can be analyzed, if the data tables are formatted according to the description in each function. The ReplMatch() function matches replicates in data sets in order to evaluate genotyping success. The GetReplTable() and GetReplStats() functions perform such an evaluation. The HpltFind() function infers putative haplotypes from families in the data set. The GetHpltTable() and GetHpltStats() functions evaluate the accuracy of the haplotype inference. The PapaDiv() function compares parent pairs in the data set and calculate their joint MHC diversity, taking into account sequence variants that occur in both parents. The CalcPdist() function calculates the p-distances from pairwise comparisons of all sequences in a data set, and mean p-distances of all pairwise comparisons within each sample in a data set. The function includes the options to specify which codons to compare and to calculate amino acid p-distances. The CreateFas() function creates a fasta file with all the sequences in the data set. The CreateSamplesFas() function creates a fasta file for each sample in the data set.

License MIT + file LICENSE

Encoding UTF-8

LazyData true

Imports rlist (>= 0.4.6.1), utils

RoxygenNote 6.1.1

NeedsCompilation no

Author Jacob Roved [aut, cre]

Maintainer Jacob Roved <jacob.roved@biol.lu.se>

Repository CRAN

Date/Publication 2019-08-11 15:20:02 UTC

R topics documented:

CalcPdist	2
CreateFas	3
CreateSamplesFas	4
GetHpltStats	5
GetHpltTable	5
GetReplStats	6
GetReplTable	7
HpltFind	7
nest_table	8
PapaDiv	9
parents_table	10
replicates_table	11
ReplMatch	11
sequence_table	12
sequence_table_fas	13
sequence_table_repl	13

Index	15
--------------	-----------

CalcPdist	<i>CalcPdist() function</i>
-----------	-----------------------------

Description

[CalcPdist](#) calculates p-distances from pairwise sequence comparisons and mean p-distances for each sample in a 'dada2' sequence table.

Usage

```
CalcPdist(seq_file, path_out, aa_pdist = NULL, codon_pos = NULL,
          input_fasta = NULL)
```

Arguments

seq_file	seq_file is a sequence table as output by the 'dada2' pipeline, which has samples in rows and nucleotide sequence variants in columns. Optionally, a fasta file can be supplied as input in the format rendered by e.g. read.fasta() from the package 'seqinr'.
path_out	is a user defined path to the folder where the output files will be saved.

aa_pdist	optional, a logical (TRUE/FALSE) that determines whether nucleotide sequences should be translated to amino acid sequences before p-distance calculation, default is NULL/FALSE.
codon_pos	optional, a vector of codon positions to include in p-distance calculations, if this argument is omitted, p-distance calculations are made using all codons.
input_fasta	optional, a logical (TRUE/FALSE) that indicates whether the input file is a fasta file (TRUE) or a dada2 sequence table (NULL/FALSE), default is NULL/FALSE.

Value

The function returns a matrix with p-distances of all pairwise sequence comparisons. This table is saved as a .csv file in the output path. If a fasta file is used as input, only the p-distance matrix will be produced. If a sequence table is given as input file, the function additionally returns a table with the mean p-distance for each sample. If a sequence table is given as input file, the sequences are named in the output matrix by an index number corresponding to their column number in the sequence table.

See Also

For more information about 'dada2' visit <<https://benjjneb.github.io/dada2>>

Examples

```
seq_file <- sequence_table_fas
path_out <- tempdir()
CalcPdist(seq_file, path_out, aa_pdist=NULL, codon_pos=c(1,2,3,4,5,6,7,8), input_fasta=NULL)
```

CreateFas	<i>CreateFas() function</i>
-----------	-----------------------------

Description

`CreateFas` creates a FASTA file with all the sequences in a 'dada2' sequence table.

Usage

```
CreateFas(seq_table, path_out)
```

Arguments

seq_table	seq_table is a sequence table as output by the 'dada2' pipeline, which has samples in rows and nucleotide sequence variants in columns.
path_out	is a user defined path to the folder where the output files will be saved.

Value

A FASTA file with all the sequences in a 'dada2' sequence table. The sequences are named in the FASTA file by an index number corresponding to their column number in the sequence table.

See Also

[CreateSamplesFas](#); for more information about 'dada2' visit <<https://benjjneb.github.io/dada2>>

Examples

```
seq_table <- sequence_table_fas
path_out <- tempdir()
CreateFas(seq_table, path_out)
```

CreateSamplesFas *CreateSamplesFas()* function

Description

[CreateSamplesFas](#) creates a set of FASTA files with the sequences present in each sample in a 'dada2' sequence table.

Usage

```
CreateSamplesFas(seq_table, path_out)
```

Arguments

seq_table	seq_table is a sequence table as output by the 'dada2' pipeline, which has samples in rows and nucleotide sequence variants in columns.
path_out	is a user defined path to the folder where the output files will be saved.

Value

A set of FASTA files with the sequences present in each sample in the sequence table. The sequences are named in the FASTA files by an index number corresponding to their column number in the sequence table, thus identical sequences will have identical sample names in all the FASTA files.

See Also

[CreateFas](#); for more information about 'dada2' visit <<https://benjjneb.github.io/dada2>>

Examples

```
seq_table <- sequence_table_fas
path_out <- tempdir()
CreateSamplesFas(seq_table, path_out)
```

GetHpltStats	<i>GetHpltStats() function</i>
--------------	--------------------------------

Description

[GetHpltStats](#) uses the output files produced by the [HpltFind\(\)](#) function to calculate the mean of the mean proportion of incongruent sequences across all nests in the data set.

Usage

```
GetHpltStats(filepath)
```

Arguments

`filepath` is a user defined path to the folder where the output files from the [HpltFind\(\)](#) function have been saved.

Value

A mean of the mean proportion of incongruent sequences for each nest.

See Also

[HpltFind](#); [GetHpltTable](#)

Examples

```
filepath <- system.file("extdata/HpltFindOut/", package="MHCtools")
GetHpltStats(filepath)
```

GetHpltTable	<i>GetHpltTable() function</i>
--------------	--------------------------------

Description

[GetHpltTable](#) uses the output files produced by the [HpltFind\(\)](#) function to produce a table with the mean proportion of incongruent sequences for each nest. If the mean proportion of incongruent sequences is generally low, but certain nests have many incongruent sequences, biological reasons may be causing the mismatches, e.g. extra-pair fertilizations or recombination events.

Usage

```
GetHpltTable(filepath)
```

Arguments

filepath is a user defined path to the folder where the output files from the HpltFind() function have been saved.

Value

A table with the mean proportion of incongruent sequences for each nest.

See Also

[HpltFind](#); [GetHpltStats](#)

Examples

```
filepath <- system.file("extdata/HpltFindOut/", package="MHCtools")
GetHpltTable(filepath)
```

GetReplStats

GetReplStats function

Description

[GetReplStats](#) uses the output files produced by the ReplMatch() function to calculate statistics on the agreement between replicated samples in the sequencing experiment.

Usage

```
GetReplStats(filepath)
```

Arguments

filepath is a user defined path to the folder where the output files from the ReplMatch() function have been saved.

Value

A list containing the number of replicate sets with zero incongruent sequences, the proportion of replicate sets with zero incongruent sequences, the mean of the mean proportion of incongruent sequences across all replicate sets, and the repeatability of the sequencing experiment.

See Also

[ReplMatch](#); [GetReplTable](#)

Examples

```
filepath <- system.file("extdata/ReplMatchOut/", package="MHCtools")
GetReplStats(filepath)
```

GetReplTable	<i>GetReplTable function</i>
--------------	------------------------------

Description

[GetReplTable](#) uses the output files produced by the `ReplMatch()` function to produce a table with the replicate sets and their respective mean proportion of incongruent sequences.

Usage

```
GetReplTable(filepath)
```

Arguments

`filepath` is a user defined path to the folder where the output files from the `ReplMatch()` function have been saved.

Value

A table with the mean proportion of incongruent sequences for each replicate set.

See Also

[ReplMatch](#); [GetReplStats](#)

Examples

```
filepath <- system.file("extdata/ReplMatchOut/", package="MHCtools")
GetReplTable(filepath)
```

HpltFind	<i>HpltFind() function</i>
----------	----------------------------

Description

[HpltFind](#) is designed to automatically infer major histocompatibility complex (MHC) haplotypes from the genotypes of parents and offspring in families (defined as nests) in non-model species, where MHC sequence variants cannot be identified as belonging to individual loci. The functions `GetHpltTable()` and `GetHpltStats()` are designed to evaluate the output files.

Usage

```
HpltFind(nest_table, seq_table, path_out)
```

Arguments

nest_table	is a table containing the sample names of parents and offspring in each nest. This table should be organized so that the individual names are in the first column (Sample_ID), and the nest number is in the second column (Nest). For each nest, the first two rows should be the parents, followed immediately by the offspring in the subsequent rows, and then followed by the next nest, and so on. It is assumed that nests are numbered consecutively beginning at 1.
seq_table	seq_table is a sequence table as output by the 'dada2' pipeline, which has samples in rows and nucleotide sequence variants in columns.
path_out	is a user defined path to the folder where the output files will be saved.

Value

A set of R lists containing for each nest the putative haplotypes, the names of sequences that could not be resolved with certainty in each parent, the names of the sequences that were incongruent in the genotypes of the nest, and the mean proportion of incongruent sequences (which is a measure of the haplotype inference success and largely influenced by the exactness of the genotyping experiment). The sequences are named in the output by an index number corresponding to their column number in the sequence table, thus identical sequences will have identical sample names in all the output files. These files can be reopened in R using e.g. the list.load() function in the 'rlist' package.

See Also

[GetHpltTable](#); [GetHpltStats](#); for more information about 'dada2' visit <<https://benjjneb.github.io/dada2>>

Examples

```
nest_table <- nest_table
seq_table <- sequence_table
path_out <- tempdir()
HpltFind(nest_table, seq_table, path_out)
```

nest_table

Data nest_table

Description

nest_table, parents_table, and sequence_table comprise a randomized dataset derived from a real major histocompatibility complex (MHC) genotyping experiment with data from parents and offspring. The nucleotide sequences have been replaced with randomly generated sequences, and sample names have been anonymized.

Usage

```
nest_table
```


Format

nest_table is a data frame with 213 samples in rows and 2 columns:

Sample_ID Sample ID
Nest Nest index number

Source

original data.

 PapaDiv

PapaDiv() function

Description

PapaDiv calculates the joint major histocompatibility complex (MHC) diversity in parent pairs, taking into account alleles that are shared between the parents. The joint diversity in parent pairs is useful for heritability analyses in non-model species, where one wants to estimate the heritability of MHC diversity. The number of unique alleles in offspring may not be directly derived from the parental genotypes if some alleles are shared between the parents.

Usage

```
PapaDiv(parents_table, seq_table, path_out)
```

Arguments

parents_table is a table containing the sample names of the parents in each nest. This table should be organized so that each row represents one nest, with the individual names of the mothers in the first column (Mother), and the individual names of the fathers in the second column (Father).

seq_table seq_table is a sequence table as output by the 'dada2' pipeline, which has samples in rows and nucleotide sequence variants in columns.

path_out is a user defined path to the folder where the output files will be saved.

Value

a set of R lists containing for the joint diversity of each parent pair, the proportion of sequences that are shared between the parents, the diversity of each of the parents, the observed sequence variants in each parent, the matched sequence variants, and the incongruent sequence variants in each parent. The sequences are named in the output by an index number corresponding to their column number in the sequence table, thus identical sequences will have identical sample names in all the output files. These files are saved in a sub folder in the output path called Parent_pairs (created by PapaDiv()) and can be reopened in R using the list.load() function in the 'rlist' package. For downstream data analysis, the PapaDiv() function also produces a summary table with the names of the parents in a pair, their respective MHC diversities, and the joint parent pair diversity. This table is saved as a .csv file in the output path.

See Also

For more information about 'dada2' visit <<https://benjjneb.github.io/dada2>>

Examples

```
parents_table <- parents_table
seq_table <- sequence_table
path_out <- tempdir()
PapaDiv(parents_table, seq_table, path_out)
```

parents_table	<i>Data parents_table</i>
---------------	---------------------------

Description

nest_table, parents_table, and sequence_table comprise a randomized dataset derived from a real major histocompatibility complex (MHC) genotyping experiment with data from parents and offspring. The nucleotide sequences have been replaced with randomly generated sequences, and sample names have been anonymized.

Usage

```
parents_table
```

Format

parents_table is a data frame with 57 parent pairs in rows and 2 columns:

Mother Mother ID

Father Father ID

Source

original data.

replicates_table	<i>Data replicates_table</i>
------------------	------------------------------

Description

replicates_table and sequence_table_repl comprise a randomized dataset derived from a real major histocompatibility complex (MHC) genotyping experiment with technical replicates. The nucleotide sequences have been replaced with randomly generated sequences, and sample names have been anonymized.

Usage

```
replicates_table
```

Format

replicates_table is a data frame with 111 technical replicate samples in rows and 2 columns:

Sample_ID Technical replicate sample ID

Replic_set Index number of replicate set

Source

original data.

ReplMatch	<i>ReplMatch() function</i>
-----------	-----------------------------

Description

In amplicon filtering it is sometimes valuable to compare technical replicates in order to estimate the accuracy of a genotyping experiment. This may be done both to optimize filtering settings and to estimate repeatability to report in a publication. [ReplMatch](#) is designed to automatically compare technical replicates in an amplicon filtering data set and report the proportion of mismatches. The functions `GetReplTable()` and `GetReplStats()` are designed to evaluate the output files.

Usage

```
ReplMatch(repl_table, seq_table, path_out)
```

Arguments

repl_table	is a table containing the sample names of technical replicates in the data set. This table should be organized so that the individual names are in the first column (Sample_ID), and the index number of the replicate set is in the second column (Replic_set). Replicate sets are allowed to contain more than two replicates. It is assumed that replicate sets are numbered consecutively beginning at 1.
seq_table	seq_table is a sequence table as output by the 'dada2' pipeline, which has samples in rows and nucleotide sequence variants in columns.
path_out	is a user defined path to the folder where the output files will be saved.

Value

A set of R lists containing for each replicate set the observed sequence variants, the names of the sequences that were incongruent in the replicates, and the mean proportion of incongruent sequences (if 100 matches are expected between the replicates, this is equivalent of an error rate in the sequencing process). The sequences are named in the output by an index number corresponding to their column number in the sequence table, thus identical sequences will have identical sample names in all the output files. These files can be reopened in R using e.g. the list.load() function in the 'rlist' package.

See Also

[GetReplTable](#); [GetReplStats](#); for more information about 'dada2' visit <<https://benjjneb.github.io/dada2>>

Examples

```
repl_table <- replicates_table
seq_table <- sequence_table_repl
path_out <- tempdir()
ReplMatch(repl_table, seq_table, path_out)
```

sequence_table	<i>Data sequence_table</i>
----------------	----------------------------

Description

nest_table, parents_table, and sequence_table comprise a randomized dataset derived from a real major histocompatibility complex (MHC) genotyping experiment with data from parents and offspring. The nucleotide sequences have been replaced with randomly generated sequences, and sample names have been anonymized.

Usage

```
sequence_table
```

Format

sequence_table is a data frame with 334 samples in rows and 329 DNA sequence variants in columns.

Source

original data.

sequence_table_fas *Data sequence_table_fas*

Description

sequence_table_fas is a randomized dataset derived from a real major histocompatibility complex (MHC) genotyping experiment. The nucleotide sequences have been replaced with randomly generated sequences, and sample names have been anonymized.

Usage

sequence_table_fas

Format

sequence_table_fas is a data frame with 100 samples in rows and 166 DNA sequence variants in columns.

Source

original data.

sequence_table_repl *Data sequence_table_repl*

Description

replicates_table and sequence_table_repl comprise a randomized dataset derived from a real major histocompatibility complex (MHC) genotyping experiment with technical replicates. The nucleotide sequences have been replaced with randomly generated sequences, and sample names have been anonymized.

Usage

sequence_table_repl

Format

`sequence_table_repl` is a data frame with 412 samples in rows and 511 DNA sequence variants in columns.

Source

original data.

Index

*Topic **datasets**

- nest_table, 8
- parents_table, 10
- replicates_table, 11
- sequence_table, 12
- sequence_table_fas, 13
- sequence_table_repl, 13

- CalcPdist, 2, 2
- CreateFas, 3, 3, 4
- CreateSamplesFas, 4, 4

- GetHpltStats, 5, 5, 6, 8
- GetHpltTable, 5, 5, 8
- GetReplStats, 6, 6, 7, 12
- GetReplTable, 6, 7, 7, 12

- HpltFind, 5–7, 7

- nest_table, 8

- PapaDiv, 9, 9
- parents_table, 10

- replicates_table, 11
- ReplMatch, 6, 7, 11, 11

- sequence_table, 12
- sequence_table_fas, 13
- sequence_table_repl, 13