

Package ‘corrDNA’

October 12, 2022

Type Package

Title Finding Associations in Position-Wise Aligned DNA Sequence Dataset

Version 1.0.1

Date 2018-03-08

Author Prabina Kumar Meher & A. R. Rao

Maintainer Prabina Kumar Meher <meherprabin@yahoo.com>

Depends R(>= 3.3.0)

Imports mvtnorm

LazyData TRUE

Description Can be useful for finding associations among different positions in a position-wise aligned sequence dataset. The approach adopted for finding associations among positions is based on the latent multivariate normal distribution.

License GPL (>= 2)

NeedsCompilation no

Repository CRAN

Date/Publication 2018-04-05 15:31:17 UTC

R topics documented:

assoc_comb	2
assoc_Zi.Zj	3
assoc_Zi.ZjR	4
assoc_Zi.ZjY	5
assoc_ZiR.ZjR	6
assoc_ZiR.ZjY	7
assoc_ZiY.ZjY	8
don_dat	9

Index	10
--------------	-----------

 assoc_comb

Complete association matrix.

Description

All the six possible association matrices can be merged in to a single matrix to visualize the overall association among positions as well as among the occurrences of nucleotides of different positions, in a position-wise aligned sequence dataset.

Usage

```
assoc_comb(x, rZiZj, rZiZjR, rZiZjY, rZiRZjR, rZiRZjY, rZiYZjY)
```

Arguments

x	A dataframe of position wise aligned sequence dataset having A, T, G and C only.
rZiZj	An object generated by using the function <code>assoc_Zi.Zj</code> .
rZiZjR	An object generated by using the function <code>assoc_Zi.ZjR</code> .
rZiZjY	An object generated by using the function <code>assoc_Zi.ZjY</code> .
rZiRZjR	An object generated by using the function <code>assoc_ZiR.ZjR</code> .
rZiRZjY	An object generated by using the function <code>assoc_ZiR.ZjY</code> .
rZiYZjY	An object generated by using the function <code>assoc_ZiY.ZjY</code> .

Details

All the six association matrices are required to be generated prior to merging them into a single matrix.

Value

A numeric matrix of order $3L$ by $3L$ for the dataset of L nucleotides long sequences.

Author(s)

Prabina Kumar Meher & A. R. Rao

Examples

```
data(don_dat)
kk <- don_dat[1:300,]
zizj <- assoc_Zi.Zj(x=kk)
zizjr <- assoc_Zi.ZjR(x=kk, rZiZj=zizj)
zizjy <- assoc_Zi.ZjY(x=kk, rZiZj=zizj)
zirzjr <- assoc_ZiR.ZjR(x=kk, rZiZj=zizj, rZiZjR=zizjr)
zirzjy <- assoc_ZiR.ZjY(x=kk, rZiZj=zizj, rZiZjR=zizjr, rZiZjY=zizjy)
```

```

zizjy <- assoc_ZiY.ZjY(x=kk,rZiZj=zizj,rZiZjY=zizjy)
fin_corr <- assoc_comb(x=kk, rZiZj=zizj,rZiZjR=zizjr,rZiZjY=zizjy,
rZiRZjR=zirzjr,rZiRZjY=zirzjy,rZiYZjY=zizjy)
fin_corr

```

assoc_Zi.Zj

Association between variable Z_i and Z_j .

Description

Finding association between variables of i^{th} position and j^{th} position. In any position wise aligned sequence dataset, occurrences of $R=(A,G)$ and $Y=(C, T)$ at each position can be explained by a standard normal variate Z based on certain threshold value. So, an association between any two position in the dataset can be obtained which will be the association between the two standard normal variate at this two positions. However, the two normal variates representing the occurrences of R and Y are independent of each other at a given position.

Usage

```
assoc_Zi.Zj(x)
```

Arguments

x A dataframe of position wise aligned sequence dataset having A, T, G and C only.

Details

The user has to supply the sequence dataset in tab delimited format and not in FASTA format. Each sequence (row) should contain only standard nucleotides (A, T, G and C). Each sequence should be same length.

Value

A numeric matrix of order L by L for the dataset of L nucleotides long sequences.

Author(s)

Prabina Kumar Meher & A. R. Rao

Examples

```

data(don_dat)
kk <- don_dat[1:300,]
zizj <- assoc_Zi.Zj(x=kk)
zizj

```

 assoc_Zi.ZjR

 Association between variable Z_i and Z_jR .

Description

Finding association between variable Z at i^{th} position and Z_R at j^{th} position. Here, the standard normal variable Z represents the occurrence of $R=(A,G)$ and $Y=(C, T)$ at each position in the position wise aligned dataset, whereas the standard normal variable Z_R represents the occurrences of nucleotides A and G at any position based on some threshold value.

Usage

```
assoc_Zi.ZjR(x, rZiZj)
```

Arguments

x A dataframe of position wise aligned sequence dataset having A, T, G and C only.

rZiZj An object generated by using the function [assoc_Zi.Zj](#).

Details

The user has to supply the input dataset as well as the output generated from the function [assoc_Zi.Zj](#).

Value

A numeric matrix of order L by L for the dataset of L nucleotides long sequences.

Note

It may happen that the convergence will not reach after a certain number of iterations and will not produce any output. In such situation, the user is advised to exclude or include some positions, or otherwise include or exclude certain sequences. The user should exploit both options till convergence is reached.

Author(s)

Prabina Kumar Meher & A. R. Rao

Examples

```
data(don_dat)
kk <- don_dat[1:300,]
zizj <- assoc_Zi.Zj(x=kk)
zizjr <- assoc_Zi.ZjR(x=kk, rZiZj=zizj)
zizjr
```

assoc_Zi.ZjY Association between variable Z_i and Z_j .

Description

Finding association between variable Z at i^{th} position and Z_Y at j^{th} position. Here, the standard normal variable Z represents the occurrence of $R=(A,G)$ and $Y=(C, T)$ at each position in the position wise aligned dataset, whereas the standard normal variable Z_R represents the occurrences of nucleotides A and G at any position based on some threshold values.

Usage

```
assoc_Zi.ZjY(x, rZiZj)
```

Arguments

x A dataframe of position wise aligned sequence dataset having A, T, G and C only.

rZiZj An object generated by using the function [assoc_Zi.Zj](#).

Details

The user has to supply the input dataset as well as the output generated from the function [assoc_Zi.Zj](#).

Value

A numeric matrix of order L by L for the dataset of L nucleotides long sequences.

Note

It may happen that the convergence will not reach after a certain number of iterations and will not produce any output. In such situation, the user is advised to exclude or include some positions, or otherwise include or exclude certain sequences. The user should exploit both options till convergence is reached.

Author(s)

Prabina Kumar Meher & A. R. Rao

Examples

```
data(don_dat)
kk <- don_dat[1:300,]
zizj <- assoc_Zi.Zj(x=kk)
zizjy <- assoc_Zi.ZjY(x=kk, rZiZj=zizj)
zizjy
```

assoc_ZiR.ZjR *Association between variable Z_{iR} and Z_{jR} .*

Description

Finding association between variable Z_R at i^{th} position and Z_R at j^{th} position. Here, the standard normal variable Z_R represents the occurrences of nucleotides A and G at any position based on some threshold value.

Usage

```
assoc_ZiR.ZjR(x, rZiZj, rZiZjR)
```

Arguments

x	A dataframe of position wise aligned sequence dataset having A, T, G and C only.
rZiZj	An object generated by using the function assoc_Zi.Zj .
rZiZjR	An object generated by using the function assoc_Zi.ZjR .

Details

The user has to supply the input dataset as well as the outputs generated from the functions [assoc_Zi.Zj](#) and [assoc_Zi.ZjR](#).

Value

A numeric matrix of order L by L for the dataset of L nucleotides long sequences.

Note

It may happen that the convergence will not reach after a certain number of iterations and will not produce any output. In such situation, the user is advised to exclude or include some positions, or otherwise include or exclude certain sequences. The user should exploit both options till convergence is reached.

Author(s)

Prabina Kumar Meher & A. R. Rao

Examples

```
data(don_dat)
kk <- don_dat[1:300,]
zizj <- assoc_Zi.Zj(x=kk)
zizjr <- assoc_Zi.ZjR(x=kk, rZiZj=zizj)
zirzjr <- assoc_ZiR.ZjR(x=kk, rZiZj=zizj, rZiZjR=zizjr)
```

zirzjr

 assoc_ZiR.ZjY *Association between variable Z_{iR} and Z_{jY} .*

Description

Finding association between variable Z_R at i^{th} position and Z_Y at j^{th} position. Here, the standard normal variable Z_Y represents the occurrences C and T at each position in the position wise aligned dataset, and the standard normal variable Z_R represents the occurrences of nucleotides A and G at any position based on some threshold values.

Usage

```
assoc_ZiR.ZjY(x, rZiZj, rZiZjR, rZiZjY)
```

Arguments

x	A dataframe of position wise aligned sequence dataset having A, T, G and C only.
rZiZj	An object generated by using the function assoc_Zi.Zj .
rZiZjR	An object generated by using the function assoc_Zi.ZjR .
rZiZjY	An object generated by using the function assoc_Zi.ZjY .

Details

The user has to supply the input dataset as well as the outputs generated from the functions [assoc_Zi.Zj](#), [assoc_Zi.ZjR](#) and [assoc_Zi.ZjY](#).

Value

A numeric matrix of order L by L for the dataset of L nucleotides long sequences.

Note

It may happen that the convergence will not reach after a certain number of iterations and will not produce any output. In such situation, the user is advised to exclude or include some positions, or otherwise include or exclude certain sequences. The user should exploit both options till convergence is reached.

Author(s)

Prabina Kumar Meher & A. R. Rao

Examples

```

data(don_dat)
kk <- don_dat[1:300,]
zizj <- assoc_Zi.Zj(x=kk)
zizjr <- assoc_Zi.ZjR(x=kk, rZiZj=zizj)
zizjy <- assoc_Zi.ZjY(x=kk, rZiZj=zizj)
zirzjy <- assoc_ZiR.ZjY(x=kk, rZiZj=zizj, rZiZjR=zizjr, rZiZjY=zizjy)
zirzjy

```

assoc_ZiY.ZjY

Association between variable Z_{iY} and Z_{jY} .

Description

Finding association between variable Z_Y at i^{th} position and Z_Y at j^{th} position. Here, the standard normal variable Z_Y represents the occurrences of nucleotides C and T at any position based on some threshold values.

Usage

```
assoc_ZiY.ZjY(x, rZiZj, rZiZjY)
```

Arguments

x	A dataframe of position wise aligned sequence dataset having A, T, G and C only.
rZiZj	An object generated by using the function assoc_Zi.Zj .
rZiZjY	An object generated by using the function assoc_Zi.ZjY .

Details

The user has to supply the input dataset as well as the outputs generated from the functions [assoc_Zi.Zj](#) and [assoc_Zi.ZjY](#).

Value

A numeric matrix of order L by L for the dataset of L nucleotides long sequences.

Note

It may happen that the convergence will not reach after a certain number of iterations and will not produce any output. In such situation, the user is advised to exclude or include some positions, or otherwise include or exclude certain sequences. The user should exploit both options till convergence is reached.

Author(s)

Prabina Kumar Meher & A. R. Rao

Examples

```
data(don_dat)
kk <- don_dat[1:300,]
zizj <- assoc_Zi.Zj(x=kk)
zizjy <- assoc_Zi.ZjY(x=kk, rZiZj=zizj)
ziyzjy <- assoc_ZiY.ZjY(x=kk, rZiZj=zizj, rZiZjY=zizjy)
ziyzjy
```

don_dat

A sample dataset of human donor splice sites.

Description

This dataset comprises 1000 donor splice site sequences, where each sequence is of length 20 with 10 at the exon end and 10 at the intron start excluding the conserved di-nucleotide GT at the beginning of intron. This dataset was randomly taken from true donor splice sites of HS3D dataset.

Usage

```
data(don_dat)
```

References

Pollastro P, Rampone S: HS3D: Homosapiens Splice Site Data Set. *Nucleic Acids Res.* 2003, Molecular Biology Database Collection entry number 36; Annual Database Issue.

Examples

```
data(don_dat)
```

Index

assoc_comb, 2
assoc_Zi.Zj, 2, 3, 4–8
assoc_Zi.ZjR, 2, 4, 6, 7
assoc_Zi.ZjY, 2, 5, 7, 8
assoc_ZiR.ZjR, 2, 6
assoc_ZiR.ZjY, 2, 7
assoc_ZiY.ZjY, 2, 8

don_dat, 9