

Package ‘cleanUpdTSeq’

January 19, 2021

Type Package

Title This package classifies putative polyadenylation sites as true or false/internally oligodT primed

Version 1.28.0

Date 2020-07-09

Author Sarah Sheppard, Jianhong Ou, Nathan Lawson, Lihua Julie Zhu

Maintainer Jianhong Ou <Jianhong.Ou@duke.edu>; Lihua Julie Zhu
<Julie.Zhu@umassmed.edu>

Depends R (>= 3.5.0), BiocGenerics (>= 0.1.0), methods, stats

Imports BSgenome, GenomicRanges, seqinr, e1071, GenomeInfoDb, IRanges, utils, BSgenome.Drerio.UCSC.danRer7

Suggests BiocStyle, knitr, RUnit

Description This package implements a Naive Bayes classifier for accurate identification of polyadenylation sites (pA sites) from oligodT based 3 prime end sequencing such as PAS-Seq, PolyA-Seq and RNA-Seq. The classifier is highly accurate and outperforms heuristic methods.

License GPL-2

biocViews Sequencing, SequenceMatching, Genetics, GeneRegulation

VignetteBuilder knitr

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

git_branch RELEASE_3_12

git_last_commit 0a193f2

git_last_commit_date 2020-10-27

Date/Publication 2021-01-18

R topics documented:

cleanUpdTSeq-package	2
BED2GRangesSeq	3
buildClassifier	4
buildFeatureVector	5
classifier	6
data.NaiveBayes	7
featureVector-class	8

getDownstreamSequence	8
getUpstreamSequence	9
modelInfo-class	10
naiveBayes-class	10
PASclassifier-class	11
predictTestSet	11

Index	14
--------------	-----------

cleanUpdTSeq-package *This package classifies putative polyadenylation sites.*

Description

3'ends of transcripts have generally been poorly annotated. With the advent of deep sequencing, many methods have been developed to identify 3'ends. The majority of these methods use an oligodT primer which can bind to internal adenine-rich sequences, and lead to artifactual identification of polyadenylation sites. Heuristic filtering methods rely on a certain number of As downstream of a putative polyadenylation site to classify the site as true or oligodT primed. This package provides a robust method to classify putative polyadenylation sites using a Naive Bayes classifier.

Details

Package: cleanUpdTSeq
 Type: Package
 Version: 1.0
 Date: 2013-07-22
 License: GPL-2

Author(s)

Sarah Sheppard, Jianhong Ou, Nathan Lawson, Lihua Julie Zhu Maintainer: Sarah Sheppard <Sarah.Sheppard@umassmed.edu>
 Jianhong Ou <Jianhong.Ou@umassmed.edu>, Lihua Julie Zhu <Julie.Zhu@umassmed.edu>

References

1. Meyer, D., et al., e1071: Misc Functions of the Department of Statistics (e1071), TU Wien. 2012.
2. Pages, H., BSgenome: Infrastructure for Biostrings-based genome data packages.
3. Sheppard, S., Lawson, N.D. and Zhu, L.J., 2013. Accurate identification of polyadenylation sites from 3' end deep sequencing using a naive Bayes classifier. *Bioinformatics*, 29(20), pp.2564-2571.

Examples

```
#read in a test set
#### first install the package using the following command
#### BiocManager::install("cleanUpdTSeq")
if (interactive())
```

```

{
library(cleanUpdTSeq)
testFile = system.file("extdata", "test.bed", package="cleanUpdTSeq")
testSet = read.table(testFile, sep = "\t", header = TRUE)

#convert the test set to GRanges with upstream and downstream sequence information
peaks = BED2GRangesSeq(testSet,upstream.seq.ind = 7, downstream.seq.ind = 8, withSeq=TRUE)
#build the feature vector for the test set with sequence information
library(BSgenome.Drerio.UCSC.danRer7)
testSet.NaiveBayes = buildFeatureVector(peaks,BSgenomeName = Drerio, upstream = 40,
  downstream = 30, wordSize = 6, alphabet=c("ACGT"),
  sampleType = "unknown",replaceNAdistance = 30,
  method = "NaiveBayes", ZeroBasedIndex = 1, fetchSeq = FALSE)

#convert the test set to GRanges without upstream and downstream sequence information
  peaks = BED2GRangesSeq(testSet,withSeq=FALSE)

#build the feature vector for the test set without sequence information
testSet.NaiveBayes = buildFeatureVector(peaks,BSgenomeName = Drerio, upstream = 40,
  downstream = 30, wordSize = 6, alphabet=c("ACGT"),
  sampleType = "unknown",replaceNAdistance = 30,
  method = "NaiveBayes", ZeroBasedIndex = 1, fetchSeq = TRUE)

#predict the test set
data(data.NaiveBayes)
predictTestSet(data.NaiveBayes$Negative, data.NaiveBayes$Positive, testSet.NaiveBayes,
outputFile = "test-predNaiveBayes.tsv", assignmentCutoff = 0.5)
}

```

BED2GRangesSeq

BED2GRangesSeq

Description

This function converts an object of data.frame from a bed file with sequence information to GRanges with sequence information.

Usage

```

BED2GRangesSeq(data.BED, header = FALSE,
  upstream.seq.ind = 7, downstream.seq.ind = 8,
  withSeq)

```

Arguments

data.BED An object of data.frame from a bed file. The data.frame should at least contains 3 required fields: `chrome`, `chromStart`, `chromend`. The fourth field for "name" is suggested for keeping track of the putative polyadenylation site from the input to the output. The sixth field for "strand" is suggested, as this will affect the classification. For this function, the bed data.frame may also contain two additional fields for the sequence upstream and downstream of the putative pA site. If these are not supplied, the sequence may be obtained when the feature vector is built. Please see <http://genome.ucsc.edu/FAQ/FAQformat.html#format1> for more information about the bed file format.

header	header = Boolean TRUE if the first row is the header FALSE if the first row is data
upstream.seq.ind	upstream.seq.ind = int, to delineate the column location containing the sequence upstream of the putative pA site
downstream.seq.ind	downstream.seq.ind = int, to delineate the column location containing the sequence downstream of the putative pA site
withSeq	TRUE indicates that there are upstream and downstream sequences in the file, FALSE indicates that there is no upstream or downstream sequences in the file

Value

Returns object of GRanges

Author(s)

Sarah Sheppard, Jianhong Ou, Nathan Lawson, Lihua Julie Zhu

Examples

```
testFile <- system.file("extdata", "test.bed", package="cleanUpdTSeq")
testSet <- read.table(testFile, sep = "\t", header = TRUE)
peaks <- BED2GRangesSeq(testSet,withSeq=TRUE)
```

buildClassifier	<i>get Naive Bayes Classifier</i>
-----------------	-----------------------------------

Description

Computes the conditional a-posterior probabilities of a categorical class variable given independent predictor variables using the Bayes rule.

Usage

```
buildClassifier(Ndata.NaiveBayes, Pdata.NaiveBayes,
               upstream=40L, downstream=30L, wordSize=6L,
               genome=Drerio, alphabet=c("ACGT"))
```

Arguments

Ndata.NaiveBayes	This is the negative training data, described further in data.NaiveBayes .
Pdata.NaiveBayes	This is the positive training data, described further in data.NaiveBayes .
upstream	This is the length of upstream sequence to use in the analysis.
downstream	This is the length of downstream sequence to use in the analysis.
wordSize	This is the size of the word to use as a feature for the upstream sequence. wordSize = 6 should always be used.
genome	Name of the genome to use to get sequence. To find out a list of available genomes, please type <code>available.genomes()</code> in R.
alphabet	These are the bases to use, for example DNA bases ACTG.

Value

An object of class "naiveBayes".

Author(s)

Jianhong Ou

See Also

[naiveBayes](#)

Examples

```
if (interactive()){
  data(data.NaiveBayes)
  classifier <- buildClassifier(data.NaiveBayes$Negative, data.NaiveBayes$Positive)
}
```

buildFeatureVector *buildFeatureVector*

Description

This function creates a data frame. Fields include peak name, upstream sequence, downstream sequence, and features to be used in classifying the putative polyadenylation site.

Usage

```
buildFeatureVector(peaks, BSgenomeName = Drerio, upstream = 50L,
  downstream = 40L, wordSize = 6L, alphabet = c("ACGT"),
  sampleType = c("TP", "TN", "unknown"), replaceNAdistance = 40L,
  method = c("NaiveBayes", "SVM"), ZeroBasedIndex = 1L, fetchSeq = FALSE)
```

Arguments

peaks	An object of GRanges that may contain the upstream and downstream sequence information. This item is created by the function BED2GRangesSeq .
BSgenomeName	Name of the genome to use to get sequence. To find out a list of available genomes, please type <code>available.genomes()</code> in R.
upstream	This is the length of upstream sequence to use in the analysis.
downstream	This is the length of downstream sequence to use in the analysis.
wordSize	This is the size of the word to use as a feature for the upstream sequence. wordSize = 6 should always be used.
alphabet	These are the bases to use, for example DNA bases ACTG.
sampleType	This is the type of sample. For example TP (true positive) or TN (true negative) for training data or unknown for test data.
replaceNAdistance	If there is no A in the downstream sequence, then use this for the average distance of As to the putative polyadenylation site.

method	This is which machine learning method to specify. For this release, method should always be set to NaiveBayes.
ZeroBasedIndex	If the coordinates are set using Zero Based indexing, set this = 1.
fetchSeq	Boolean, for getting upstream and downstream sequence at this step or not.

Value

An object of "[featureVector](#)"

Author(s)

Sarah Sheppard, Jianhong Ou, Nathan Lawson, Lihua Julie Zhu

Examples

```
testFile = system.file("extdata", "test.bed", package="cleanUpdTSeq")
testSet = read.table(testFile, sep = "\t", header = TRUE)

#convert the test set to GRanges with upstream and downstream sequence information
peaks = BED2GRangesSeq(testSet[1:10, ], upstream.seq.ind = 7, downstream.seq.ind = 8, withSeq=TRUE)
#build the feature vector for the test set with sequence information
library(BSgenome.Drerio.UCSC.danRer7)
testSet.NaiveBayes = buildFeatureVector(peaks,BSgenomeName = Drerio, upstream = 40,
  downstream = 30, wordSize = 6, alphabet=c("ACGT"),
  sampleType = "unknown",replaceNAdistance = 30,
  method = "NaiveBayes", ZeroBasedIndex = 1, fetchSeq = FALSE)

#convert the test set to GRanges without upstream and downstream sequence information
peaks = BED2GRangesSeq(testSet[1:10, ],withSeq=FALSE)

#build the feature vector for the test set without sequence information
testSet.NaiveBayes = buildFeatureVector(peaks,BSgenomeName = Drerio, upstream = 40,
  downstream = 30, wordSize = 6, alphabet=c("ACGT"),
  sampleType = "unknown",replaceNAdistance = 30,
  method = "NaiveBayes", ZeroBasedIndex = 1, fetchSeq = TRUE)
```

classifier

An object of class "naiveBayes" generated from data.NaiveBayes

Description

An object of class "naiveBayes" generated from data.NaiveBayes

Usage

```
data("classifier")
```

Format

An object of class "[PASclassifier](#)" including components:

Details

```

classifier is generated by
library(BSgenome.Drerio.UCSC.danRer7)
data(data.NaiveBayes)
classifier <- buildClassifier(data.NaiveBayes$Negative, data.NaiveBayes$Positive)

```

Examples

```

data(classifier)
names(classifier)

```

data.NaiveBayes	<i>Training Data</i>
-----------------	----------------------

Description

This is the negative and positive training data.

Usage

```
data(data.NaiveBayes)
```

Format

A list with 2 data frame, "Negative" and "Positive". Negative has 9219 observations on the following 4120 variables. And Positive is a data frame with 22770 observations on the following 4120 variables. The format is:

Negative 'data.frame': 9219 obs. of 4120 variables:

Positive 'data.frame': 22770 obs. of 4120 variables:

Both of them have same structure.

y a numeric vector

n.A.Downstream a numeric vector

n.C.Downstream a numeric vector

n.T.Downstream a numeric vector

n.G.Downstream a numeric vector

avg.distanceA2PeakEnd a numeric vector

dimer: **such as AA, AC, AG, AT, CA, ... etc.** a numeric vector

heximer: **such as AAAAAA, ACGTAC, ... etc.** a factor with levels 0 1

upstream.seq a vector of sequence string

downstream.seq a vector of sequence string

Examples

```

data(data.NaiveBayes)
head(str(data.NaiveBayes$Negative))
head(str(data.NaiveBayes$Positive))

```

featureVector-class *Class "featureVector"*

Description

An object of class "featureVector" represents the output of [buildFeatureVector](#)

Objects from the Class

Objects can be created by calls of the form `new("featureVector", data, info)`.

Slots

data An object of data frame. Fields include peak name, upstream sequence, downstream sequence, and features to be used in classifying the putative polyadenylation site

info Object of class [modelInfo](#)

Methods

`$, $<-` Get or set the slot of [featureVector](#)

getDownstreamSequence *getDownstreamSequence*

Description

This function gets the sequence upstream of a putative pA site (including the site)

Usage

```
getDownstreamSequence(peaks, downstream = 20, genome)
```

Arguments

peaks	GRanges containing putative pA sites
downstream	downstream = int. This is the length of the sequence to get.
genome	BSgenomeName

Value

Returns an object of GRanges with downstream sequences.

Author(s)

Sarah Sheppard, Jianhong Ou, Nathan Lawson, Lihua Julie Zhu

Examples

```
library(BSgenome.Drerio.UCSC.danRer7)
testFile <- system.file("extdata", "test.bed", package="cleanUpdTSeq")
testSet <- read.table(testFile, sep="\t", header=TRUE)
peaks <- BED2GRangesSeq(testSet[1:10, ], withSeq=FALSE)
seq = getDownstreamSequence(peaks, downstream=30, genome=Drerio)
```

getUpstreamSequence *Get upstream sequences of the putative pA site*

Description

This function gets the sequence upstream of a putative pA site (including the site)

Usage

```
getUpstreamSequence(peaks, upstream = 40, genome)
```

Arguments

peaks	GRanges containing putative pA sites
upstream	upstream = int. This is the length of the sequence to get.
genome	BSgenomeName

Value

Returns an object of GRanges with upstream sequences.

Author(s)

Sarah Sheppard, Jianhong Ou, Nathan Lawson, Lihua Julie Zhu

Examples

```
library(BSgenome.Drerio.UCSC.danRer7)
testFile <- system.file("extdata", "test.bed", package="cleanUpdTSeq")
testSet <- read.table(testFile, sep="\t", header=TRUE)
peaks <- BED2GRangesSeq(testSet[1:10, ], withSeq=FALSE)
seq = getUpstreamSequence(peaks, upstream=40, genome=Drerio)
```

modelInfo-class	Class "modelInfo"
-----------------	-------------------

Description

An object of class "modelInfo" represents the information of sequence to use in the analysis

Objects from the Class

Objects can be created by calls of the form `new("modelInfo", upstream, downstream, wordSize, alphabe, genome)`.

Slots

genome Name of the genome to use to get sequence. To find out a list of available genomes, please type `available.genomes()` in R.

upstream This is the length of upstream sequence to use in the analysis.

downstream This is the length of downstream sequence to use in the analysis.

wordSize This is the size of the word to use as a feature for the upstream sequence. `wordSize = 6` should always be used.

alphabet These are the bases to use, for example DNA bases ACTG.

Methods

`$, $<-` Get or set the slot of `modelInfo`

naiveBayes-class	Class "naiveBayes"
------------------	--------------------

Description

An object of class "naiveBayes" represents the conditional a-posterior probabilities of a categorical class variable given independent predictor variables using the Bayes rule.

Objects from the Class

Objects can be created by calls of the form `new("naiveBayes", apriori, tables, levels, call)`.

Slots

apriori Class distribution for the dependent variable.

tables A list of tables, one for each predictor variable. For each categorical variable a table giving, for each attribute level, the conditional probabilities given the target class. For each numeric variable, a table giving, for each target class, mean and standard deviation of the (sub-)variable.

Methods

`$, $<-` Get or set the slot of `naiveBayes`

PASclassifier-class *Class "PASclassifier"*

Description

An object of class "PASclassifier" represents the output of `buildClassifier`

Objects from the Class

Objects can be created by calls of the form `new("PASclassifier", classifier, info)`.

Slots

classifier Object of class "naiveBayes" The output of `naiveBayes`. An object of class "naive-Bayes" including components:

a priori Class distribution for the dependent variable.

tables A list of tables, one for each predictor variable. For each categorical variable a table giving, for each attribute level, the conditional probabilities given the target class. For each numeric variable, a table giving, for each target class, mean and standard deviation of the (sub-)variable.

info Object of class `modelInfo`

Methods

`$, $<-` Get or set the slot of `PASclassifier`

Examples

```
data(classifier)
classifier$info$upstream
classifier$info$wordSize
classifier$info$alphabet
```

`predictTestSet` *predictTestSet*

Description

This function can be used to predict the probabilities for a set of putative pA sites.

Usage

```
predictTestSet(Ndata.NaiveBayes, Pdata.NaiveBayes, testSet.NaiveBayes, classifier=NULL,
  outputFile = "test-predNaiveBayes.tsv", assignmentCutoff = 0.5)
```

Arguments

Ndata.NaiveBayes	This is the negative training data, described further in data.NaiveBayes .
Pdata.NaiveBayes	This is the positive training data, described further in data.NaiveBayes .
classifier	An object of class PASClassifier .
testSet.NaiveBayes	This is the test data, a feature vector that has been built for Naive Bayes analysis using the function "buildFeatureVector".
outputFile	This is the name of the file the output will be written to.
assignmentCutoff	This is the cutoff used to assign whether a putative pA is true or false. This can be any floating point number between 0 and 1. For example, assignmentCutoff = 0.5 will assign an putative pA site with prob.1 > 0.5 to the True class (1), and any putative pA site with prob.1 <= 0.5 as False (0).

Value

The output is written to a tab separated file containing fields for peak name, the probability of the putative pA site being false (prob.0), the probability of the putative pA site being true (prob.1), the predicted class (0/False or 1/True) depending on the assignment cutoff, and the upstream and downstream sequence used in assessing the putative pA site.

PeakName	This is the name of the putative pA site (originally from the 4th field in the bed file).
prob False/oligodT internally primed	This is the probability that the putative pA site is false. Values range from 0-1, with 1 meaning the site is False/oligodT internally primed.
prob True	This is the probability that the putative pA site is true. Values range from 0-1, with 1 meaning the site is True.
pred.class	This is the predicted class of the putative pA site, based on the assignment cutoff. 0= False/oligodT internally primed, 1 = True
UpstreamSeq	This is the upstream sequence of the putative pA site used in the analysis.
DownstreamSeq	This is the downstream sequence of the putative pA site used in the analysis.

The function also return an invisible matrix including all info as decribed above.

Author(s)

Sarah Sheppard, Jianhong Ou, Nathan Lawson, Lihua Julie Zhu

References

Sarah Sheppard, Nathan D. Lawson, and Lihua Julie Zhu. 2013. Accurate identification of polyadenylation sites from 3' end deep sequencing using a naive Bayes classifier. *Bioinformatics*. Under revision

Examples

```
testFile = system.file("extdata", "test.bed", package="cleanUpdTSeq")
testSet = read.table(testFile, sep = "\t", header = TRUE)

#convert the test set to GRanges without upstream and downstream sequence information
peaks = BED2GRangesSeq(testSet,withSeq=FALSE)

#build the feature vector for the test set without sequence information
library(BSgenome.Drerio.UCSC.danRer7)
testSet.NaiveBayes = buildFeatureVector(peaks,BSgenomeName = Drerio, upstream = 40,
    downstream = 30, wordSize = 6, alphabet=c("ACGT"),
    sampleType = "unknown",replaceNAdistance = 30,
    method = "NaiveBayes", ZeroBasedIndex = 1, fetchSeq = TRUE)

data(data.NaiveBayes)

## sample the test data for code testing, DO NOT do this for real data
## START SAMPLING
samp <- c(1:22, sample(23:4119, 50), 4119, 4120)
Ndata.NaiveBayes <- data.NaiveBayes$Negative[,samp]
Pdata.NaiveBayes <- data.NaiveBayes$Positive[,samp]
testSet.NaiveBayes@data <- testSet.NaiveBayes@data[, samp-1]
## END SAMPLING

predictTestSet(Ndata.NaiveBayes,
    Pdata.NaiveBayes,
    testSet.NaiveBayes,
    outputFile="test-predNaiveBayes.xls",
    assignmentCutoff = 0.5)
```

Index

- * **classes**
 - featureVector-class, 8
 - modelInfo-class, 10
 - naiveBayes-class, 10
 - PASClassifier-class, 11
 - * **datasets**
 - classifier, 6
 - data.NaiveBayes, 7
 - * **misc**
 - BED2GRangesSeq, 3
 - buildClassifier, 4
 - buildFeatureVector, 5
 - getDownstreamSequence, 8
 - getUpstreamSequence, 9
 - predictTestSet, 11
 - * **package**
 - cleanUpdTSeq-package, 2
 - featureVector (featureVector-class), 8
 - featureVector-class, 8
 - getDownstreamSequence, 8
 - getUpstreamSequence, 9
 - modelInfo, 8, 10, 11
 - modelInfo (modelInfo-class), 10
 - modelInfo-class, 10
 - naiveBayes, 5, 10, 11
 - naiveBayes (naiveBayes-class), 10
 - naiveBayes-class, 10
 - PASClassifier, 6, 11, 12
 - PASClassifier (PASClassifier-class), 11
 - PASClassifier-class, 11
 - predictTestSet, 11
 - featureVector (featureVector-class), 8
 - featureVector-class, 8
 - getDownstreamSequence, 8
 - getUpstreamSequence, 9
 - modelInfo, 8, 10, 11
 - modelInfo (modelInfo-class), 10
 - modelInfo-class, 10
 - naiveBayes, 5, 10, 11
 - naiveBayes (naiveBayes-class), 10
 - naiveBayes-class, 10
 - PASClassifier, 6, 11, 12
 - PASClassifier (PASClassifier-class), 11
 - PASClassifier-class, 11
 - predictTestSet, 11
- BED2GRangesSeq, 3, 5
- buildClassifier, 4, 11
- buildFeatureVector, 5, 8
- classifier, 6
- cleanUpdTSeq (cleanUpdTSeq-package), 2
- cleanUpdTSeq-package, 2
- data.NaiveBayes, 4, 7, 12
- featureVector, 6, 8