

rGADEM

October 5, 2010

GADEM

Motif Analysis with rGADEM

Description

It is an R implementation of GADEM, a powerful computational tools for de novo motif discovery.

Usage

```
gadem<-GADEM(Sequences, seed=1, genome=NULL, verbose=TRUE, numWordGroup=3, numTop3mer,
populationSize=100, pValue=0.0002, eValue=0.0, extTrim=1, minSpaceWidth=0, maxSpaceWi
```

Arguments

<code>Sequences</code>	Sequences from BED or FASTA file are converted into XString object view
<code>seed</code>	When a seed is specified, the run results are deterministic
<code>genome</code>	Specify the genome
<code>verbose</code>	Print immediate results on screen [TRUE=yes (default), FALSE=no]. These results include the motif consensus sequence, number of sites (in sequences subjected to EM optimization, see <code>-fEM</code> , above), and $\ln(E\text{-value})$.
<code>numWordGroup</code>	number of non-zero k-mer groups
<code>numTop3mer</code>	Number of top-ranked trimers for spaced dyads (default: 20).
<code>numTop4mer</code>	Number of top-ranked tetramers for spaced dyads (default: 40).
<code>numTop5mer</code>	Number of top-ranked pentamers for spaced dyads (default: 60).
<code>numGeneration</code>	Number of genetic algorithm (GA) generations (default: 5).
<code>populationSize</code>	GA population size (default: 100). Both default settings should work well for most datasets (ChIP-chip and ChIP-seq). The above two arguments are ignored in a seeded analysis, because spaced dyads and GA are no longer needed (<code>-gen</code> is set to 1 and <code>-pop</code> is set to 10 internally, corresponding to the 10 maxp choices).
<code>pValue</code>	P-value cutoff for declaring BINDING SITES (default: 0.0002). Depending on data size and the motif, you might want to assess more than one value. For ChIP-seq data (e.g., 10 thousand +/-200-bp max-center peak cores), $p=0.0002$ often seems appropriate. However, short motifs may require a less stringent setting.

eValue	ln(E-value) cutoff for selecting MOTIFS (default: 0.0). If a seeded analysis fails to identify the expected motif, run GADEM with <code>-verbose 1</code> to show motif ln(E-value)s on screen, then rerun with a larger ln(E-value) cutoff. This can help in identifying short and/or low abundance motifs, for which the default E-value threshold may be too low.
extTrim	Base extension and trimming (1 -yes, 0 -no) (default: 1).
minSpaceWidth	Minimal number of unspecified nucleotides in spaced dyads (default: 0).
maxSpaceWidth	Maximal number of unspecified nucleotides in spaced dyads (default: 10). <code>-mingap</code> and <code>-maxgap</code> control the lengths of spaced dyads, and, with <code>-extrim</code> , control motif lengths. Longer motifs can be discovered by setting <code>-maxgap</code> to larger values (e.g. 50).
useChIPscore	Use top-scoring sequences for deriving PWMs. Sequence (quality) scores are stored in sequence header (see documentation). 0 - no (default, randomly select sequences), 1 - yes.
numEM	Number of EM steps (default: 40). One might want to set it to a larger value (e.g. 80) in a seeded run, because such runs are fast.
fEM	Fraction of sequences used in EM to obtain PWMs in an unseeded analysis (default: 0.5). For unseeded motif discovery in a large dataset (e.g. >10 million nt), one might want to set <code>-fEM</code> to a smaller value (e.g., 0.3 or 0.4) to reduce run time.
widthWt	For <code>-posWt 1</code> or <code>3</code> , width of central sequence region with large EM weights for PWM optimization (default: 50). This argument is ignored when <code>-posWt</code> is 0 (uniform prior) or 2 (Gaussian prior).
fullScan	GADEM keeps two copies of the input sequences internally: one (D) for discovering PWMs and one (S) for scanning for binding sites using the PWMs. Once a motif is identified, its instances in set D are always masked by Ns. However, masking motif instances in set S is optional, and scanning unmasked sequences allows sites of discovered motifs to overlap.
userBackgModel	To run analysis in background (default : 0).
slideWinPwm	sliding window for comparing pwm similarity (default : 6).
stopCriterion	Stop analysis.
MarkovOrder	Background Markov order,user-specified order highest order available in user-specified background indicator (default : 0).
userMarkovOrder	Background Markov order,user-specified order highest order available in user-specified background indicator (default : 0).
numBackgSets	Number of sets of background sequences (default: 10). The background sequences are simulated using the [a,c,g,t] frequencies in the input sequences, with length matched between the two sets. The background sequences are used as the random sequences for assessing motif enrichment in the input data. Another set (same default: 10) of background sequences is independently generated to approximate the empirical llr score distribution when <code>-pgf</code> is set to 0.
weightType	Weight profile for positions on the sequence. 0 - no weight (uniform spatial prior, default), 1 - small or zero weights for the ends and large weights for the center (e.g. the center 50 bp). If you expect strong central enrichment (as in ChIP-seq) and your sequences are long (e.g. >200 bp), choose type 1.

pgf	By default, GADEM uses the Staden probability generating function (pgf) method to approximate the exact llr score null distribution.
startPWMfound	Value for the PWM (default : 0).
bOrder	The order of the background Markov model for computing llr scores: 0 - 0th 1 - 1st 2 - 2nd 8 - 8th
bFileName	Reading user-specified background models.
Spwm	File name for the seed PWM, when a seeded approach is used. can be used as the starting PWM for the EM algorithm. This will help find an expected motif and is much faster than unseeded de novo discovery. Also, when a seed PWM is specified, the run results are deterministic, so only a single run is needed (repeat runs with the same settings will give identical results). In contrast, unseeded runs are stochastic, and we recommend comparing results from several repeat runs.

Author(s)

Arnaud Droit <arnaud.droit@ircm.qc.ca>

Examples

```
library(BSgenome.Hsapiens.UCSC.hg18)
pwd<-"#INPUT FILES- BedFiles, FASTA, etc."
path<- system.file("extdata/Test_100.bed", package="rGADEM")
BedFile<-paste(pwd, path, sep="")
BED<-read.table(BedFile, header=FALSE, sep="\t")
BED<-data.frame(chr=as.factor(BED[,1]), start=as.numeric(BED[,2]), end=as.numeric(BED[,3]))
#Create RD files
rgBED<-IRanges(start=BED[,2], end=BED[,3])
Sequences<-RangedData(rgBED, space=BED[,1])

gadem<-GADEM(Sequences, verbose=1, genome=Hsapiens)
```

align-class

Class "align"

Description

This object contains the individual motifs identified but also the location (seqID and position) of the sites in the original sequence data. It also included the spaced dyad from which the motifs is derived, PWM score p-value cutoff for the run.

Objects from the Class

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

Slots

seq :Motif identified .
chr :Chromosome identified.
start :Sequence start.
end :Sequence end.
strand :Strand position.
seqID :Sequence identification.
pos :Position identification.
pval :p-Value for each identification.
fastaHeader :Fasta accession.

Author(s)

Arnaud Droit <arnaud.droit@ircm.qc.ca>

See Also

[gadem](#), [motif](#), [parameters](#)

Examples

```
showClass("align")
```

gadem-class

Class "gadem"

Description

This object contains all gadem output information.

Objects from the Class

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

Slots

motifList List of input PWM.
parameters List of rGADEM parameters.

Methods

[signature(x = "gadem"): subset gadem object.
[[] signature(x = "gadem"): subset gadem object.
nMotifs signature(x = "gadem"): Number of motifs identified
names signature(x = "gadem"): Assign motifs names.
dim signature(x = "gadem"): Number of sequences identified for each motifs.
consensus signature(x = "gadem"): Sequence of consensus motifs.

nOccurrences signature(x = "gadem"):View of PWMs.
startPos signature(x = "gadem"):Start position for each sequences.
endPos signature(x = "gadem"):End position for each sequences.
getPWM signature(x = "gadem"):End position for each sequences.

Author(s)

Arnaud Droit <arnaud.droit@ircm.qc.ca>

See Also

[motif](#), [align](#), [parameters](#)

Examples

```
showClass("gadem")
```

motif-class

Class "motif"

Description

This object contains contains PWM, motif consensus, motif length and all aligned sequences for a specific motif

Objects from the Class

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

Slots

pwm :PWM results.
consensus :Sequences consensus.
alignList :List of sequences alignment.
name :Name of sequences.

Author(s)

Arnaud Droit <arnaud.droit@ircm.qc.ca>

See Also

[gadem](#), [align](#), [parameters](#)

Examples

```
showClass("gadem")
```

parameters-class *Class "parameters"*

Description

This object contains contains parameters of GADEM analysis

Objects from the Class

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

Slots

numWordGroup :number of non-zero k-mer groups.

numTop3mer :Number of top-ranked trimers for spaced dyads (default: 20).

verbose :Print immediate results on screen [1=yes (default), 0=no].

numTop4mer :Number of top-ranked tetramers for spaced dyads (default: 40).

numTop5mer :Number of top-ranked pentamers for spaced dyads (default: 60).

numGeneration :Number of genetic algorithm (GA) generations (default: 5).

populationSize :GA population size (default: 100).

pValue :P-value cutoff for declaring BINDING SITES (default: 0.0002).

eValue :ln(E-value) cutoff for selecting MOTIFS (default: 0.0).

extTrim :Base extension and trimming (1 -yes, 0 -no) (default: 1).

minSpaceWidth :Minimal number of unspecified nucleotides in spaced dyads (default: 0).

maxSpaceWidth :Maximal number of unspecified nucleotides in spaced dyads (default: 10).

useChIPscore :Use top-scoring sequences for deriving PWMs.

numEM :Number of EM steps (default: 40).

fEM :Fraction of sequences used in EM to obtain PWMs in an unseeded analysis (default: 0.5).

widthWt :For -posWt 1 or 3, width of central sequence region with large EM weights for PWM optimization (default: 50).

fullScan :GADEM keeps two copies of the input sequences internally.

userBackgModel :To run analysis in background (default : 0).

slideWinPWM :sliding window for comparing pwm similarity (default : 6).

stopCriterion

MarkovOrder :Background Markov order,user-specified order highest order available in user-specified background indicator (default : 0).

userMarkovOrder :Background Markov order,user-specified order highest order available in user-specified background indicator (default : 0).

numBackgSets :Number of sets of background sequences (default: 10).

weightType :Weight profile for positions on the sequence.

pgf :By default, GADEM uses the Staden probability generating function (pgf) method to approximate the exact llr score null distribution.

startPWMfound :Value for the PWM (default : 0).

bOrder :The order of the background Markov model for computing llr scores: 0 - 0th 1 - 1st 2 - 2nd 8 - 8th

bFileName :Reading user-specified background models.

fpwm0 :File name for the seed PWM, when a seeded approach is used.

nSequences :number of input sequences.

Author(s)

Arnaud Droit <arnaud.droit@ircm.qc.ca>

See Also

[gadem](#), [align](#), [motif](#)

Examples

```
showClass("parameters")
```

Index

*Topic **GADEM**
GADEM, 1

*Topic **MOTIFS**
GADEM, 1

*Topic **classes**
align-class, 3
gadem-class, 4
motif-class, 5
parameters-class, 6

[, gadem-method (*gadem-class*), 4
[[, gadem-method (*gadem-class*), 4

align, 5, 7
align (*align-class*), 3
align-class, 3

consensus (*gadem-class*), 4
consensus, gadem-method
 (*gadem-class*), 4

dim, gadem-method (*gadem-class*), 4

endPos (*gadem-class*), 4
endPos, gadem-method
 (*gadem-class*), 4

GADEM, 1
gadem, 4, 5, 7
gadem (*gadem-class*), 4
gadem-class, 4
getPWM (*gadem-class*), 4
getPWM, gadem-method
 (*gadem-class*), 4
getPWM, motif-method
 (*gadem-class*), 4

length, gadem-method
 (*gadem-class*), 4

motif, 4, 5, 7
motif (*motif-class*), 5
motif-class, 5

names, gadem-method (*gadem-class*),
 4

names<-, gadem-method
 (*gadem-class*), 4

nMotifs (*motif-class*), 5
nMotifs, gadem-method
 (*gadem-class*), 4

nOccurrences (*gadem-class*), 4
nOccurrences, gadem-method
 (*gadem-class*), 4

parameters, 4, 5
parameters (*parameters-class*), 6
parameters, gadem-method
 (*gadem-class*), 4

parameters-class, 6
plot (*gadem-class*), 4
plot, gadem-method (*gadem-class*), 4
plot, motif-method (*gadem-class*), 4

startPos (*gadem-class*), 4
startPos, gadem-method
 (*gadem-class*), 4

summary, list-method
 (*gadem-class*), 4

viewPWM (*gadem-class*), 4
viewPWM, gadem-method
 (*gadem-class*), 4