

Local Score Package

Sébastien Déjean - Sabine Mercier - Sebastian Simon - David Robelin

2021-02-22

Contents

Introduction	2
Local Score computation methods	2
A first example: function “localScoreC()”	2
Example with real scores: function “localScoreC_double()”	4
Example using a scoring function	4
<i>p</i>-Value computation methods	5
Simulating computation: functions “monteCarlo()” and “monteCarlo_double()”	5
A mixed method: functions “karlinMonteCarlo()” and “karlinMonteCarlo_double()”	7
Exact method for integer scores: function “daudin()”	9
How to use the exact method for real scores	9
Approximate method of Karlin <i>et al.</i> : function “karlin()”	12
An improved approximate method: function “mcc()”	12
An automatic method: function “automatic_analysis()”	13
Other Functions	15
Lindley Process: to visualize optimal and suboptimal segments	15
Score Loading Function	16
Empirical distribution: function “scoreSequences2probabilityVector()”	16
Case study	16
Medium sequence	16
Short sequence	20
Large sequence	22
Several sequences	24
A larger example with a SCOP data base	25
File Formats	33
Sequence Files	33
Score Files	34
Transition Matrix Files	34
A word for Markovian model	34
library(localScore)	

Introduction

This package provides functionalities for two main tasks: 1- Calculating the local Score of a given score sequence or, of a given component sequence and a given scoring scheme. 2- Calculating the statistical relevance (p -value) of a given local Score, associated to a given sequence length and a given distribution for the model.

This first version of the package deals with a model of independent and identically distributed (I.I.D.) sequences.

Local Score computation methods

First defined in *Karlin and Altschul 1990*, it represents the value of the highest scoring segment in a sequence of scores. It corresponds to the highest cumulated sum of all subsequences (independent of length). For the score to be relevant, the expectation of the sequence should be negative. Thus, for a sequence of interest, the possible score of sequence components can be positive or negative.

A first example: function “localScoreC()”

Let us assume a score function taking its values in $[-2, -1, 0, 1, 2]$. A sample score sequence of length 100 could be

```
library(localScore)
help(localScore)
ls(pos=2)
#> [1] "CharSequence2ScoreSequence"      "CharSequences2ScoreSequences"
#> [3] "LongSeq"                          "MidSeq"
#> [5] "MySeqList"                        "RealScores2IntegerScores"
#> [7] "ShortSeq"                         "automatic_analysis"
#> [9] "daudin"                           "dico"
#> [11] "exact_mc"                         "karlin"
#> [13] "karlinMonteCarlo"                 "karlinMonteCarlo_double"
#> [15] "lindley"                           "loadMatrixFromFile"
#> [17] "loadScoreFromFile"                "localScoreC"
#> [19] "localScoreC_double"               "mcc"
#> [21] "monteCarlo"                       "monteCarlo_double"
#> [23] "scoreDictionnary2probabilityVector" "scoreSequences2probabilityVector"
#> [25] "sequences2transmatrix"             "stationary_distribution"
#> [27] "transmatrix2sequence"
sequence <- sample(-2:2, 100, replace = TRUE, prob = c(0.5, 0.3, 0.05, 0.1, 0.05))
sequence
#> [1] -1  1 -2 -1 -2 -2  0 -2 -2 -2 -2  1 -1 -1 -1 -2 -1 -2 -2 -2 -2 -2 -1  1
#> [26] -2 -2 -2 -1 -1  2 -1 -1 -2 -2 -1  0 -2  2 -2 -2  1 -1 -2  1 -2 -2 -2 -2  2
#> [51] -1 -2 -2 -2 -2 -1 -2 -2 -1 -2 -2 -2 -1 -1  1 -2 -2 -2 -1  0 -2 -1 -2 -1 -2
#> [76] -1 -2  1  2 -1 -1 -2  2 -2 -1 -2 -2  1 -1  2 -2 -1  1 -2 -1  1 -1 -1 -2 -2
localScoreC(sequence)
#> $localScore
#> value begin end
#>      3     78    79
#>
#> $suboptimalSegmentScores
#>      value begin end
```

```

#> [1,] 1 2 2
#> [2,] 1 12 12
#> [3,] 1 25 25
#> [4,] 2 31 31
#> [5,] 2 39 39
#> [6,] 1 42 42
#> [7,] 1 45 45
#> [8,] 2 50 50
#> [9,] 1 65 65
#> [10,] 3 78 79
#> [11,] 2 83 83
#> [12,] 1 88 88
#> [13,] 2 90 90
#> [14,] 1 93 93
#> [15,] 1 96 96
#>
#> $RecordTime
#> [1] 2 12 25 31 39 42 45 50 65 78 83 88 90 93 96

```

The result is a maximum score and for which subsequence it has been found: the starting position "[" and the end of it ("]"). It also yields all other subsequences with a score equal or less and its position in the \$suboptimalSegmentScores matrix. "Stopping Times" are local minima in the cumulated sum of the sequence and correspond to the beginning of excursions (potential segments of interest).

Another example with missing score values.

```

library(localScore)
sequence <- sample(c(-3,2,0,1,5), 100, replace = TRUE, prob = c(0.5, 0.3, 0.05, 0.1, 0.05))
head(sequence)
#> [1] 2 -3 -3 0 5 2
localScoreC(sequence)
#> $localScore
#> value begin end
#> 15 89 97
#>
#> $suboptimalSegmentScores
#> value begin end
#> [1,] 2 1 1
#> [2,] 13 5 10
#> [3,] 1 19 19
#> [4,] 4 23 24
#> [5,] 1 30 30
#> [6,] 3 32 34
#> [7,] 1 36 36
#> [8,] 2 38 38
#> [9,] 2 40 40
#> [10,] 2 44 44
#> [11,] 1 47 47
#> [12,] 4 49 50
#> [13,] 10 56 72
#> [14,] 2 80 80
#> [15,] 2 84 84
#> [16,] 2 86 86
#> [17,] 15 89 97
#>

```

```
#> $RecordTime
#> [1] 1 5 19 23 30 32 36 38 40 44 47 49 56 80 84 86 89
```

Example with real scores: function “localScoreC_double()”

Real scores can also be considered with a dedicated function ‘localScoreC_double’.

```
score_reels=c(-1,-0.5,0,0.5,1)
proba_score_reels=c(0.2,0.3,0.1,0.2,0.2)
sample_from_model <- function(score.sple,proba.sple, length.sple){sample(score.sple,
                                size=length.sple, prob=proba.sple, replace=TRUE)}
seq.essai=sample_from_model(score.sple=score_reels,proba.sple=proba_score_reels, length.sple=10)
localScoreC_double(seq.essai)
#> $localScore
#> value begin end
#> 2 8 10
#>
#> $suboptimalSegmentScores
#> value begin end
#> [1,] 1 4 5
#> [2,] 2 8 10
#>
#> $RecordTime
#> [1] 4 8
```

Example using a scoring function

```
# Loading a fasta protein
data(MidSeq)
MidSeq
#> [1] "MSGLSGPPARRGPFPLALLLFLLGPRVLVAISFHLPIINSRKCLREEIHKDLLVTGAYEISDQSGGAGGLRSHLKITDSAGHILYSKEDATKGKF"
# or using your own fasta sequence
#MySeqAA_P49755 <- as.character(read.table(file="P49755.fasta",skip=1)[,1])
#MySeqAA_P49755
# Loading a scoring function
data(dico)
? dico
# or using your own scoring function
# dicoKyte<- loadScoreFromFile("Kyte1982.txt")
# Transforming the amino acid sequence into a score sequence with the score function in dico file
SeqScore_P49755 <- CharSequence2ScoreSequence(MidSeq,dico)
head(SeqScore_P49755)
#> [1] 2 -1 0 4 -1 0
length(SeqScore_P49755)
#> [1] 219
# Computing the local score
localScoreC(SeqScore_P49755)
#> $localScore
#> value begin end
#> 52 14 38
#>
```

```

#> $suboptimalSegmentScores
#>      value begin end
#> [1,]      5      1  4
#> [2,]      2      9  9
#> [3,]     52     14 38
#> [4,]     17    121 124
#> [5,]      7    138 139
#> [6,]      4    144 144
#> [7,]      4    147 147
#> [8,]      4    149 149
#> [9,]      4    151 151
#> [10,]     4    154 154
#> [11,]     4    157 157
#> [12,]      9    161 162
#> [13,]      2    175 175
#> [14,]     43    186 208
#>
#> $RecordTime
#> [1]  1  9 14 121 138 144 147 149 151 154 157 161 175 186

```

File dico is read using `read.csv()` with a default option `header=TRUE`. It is so necessary that the file has one. See Section “File format” for more details.

p-Value computation methods

There are different methods available to establish the statistical significance, also called *p*-value, of the local score depending on the length and the score expectation. This value describes the probability to encounter a given local score for a given score distribution and a given sequence length. Therefore it allows to determinate if the local score in question is significant or could have been obtained by chance.

For an Identically and Independently Distributed Variables model (I.I.D.), the following methods for the calculus of the *p*-value are provided:

Simulating computation: functions “`monteCarlo()`” and “`monteCarlo_double()`”

The function `monteCarlo()` simulates a number of score sequences similar (same distribution) to the one having yielded the local score in question. Therefore, it requires a function that produces such sequences and the parameters used by this function. See the help page and the following example to use the empirical distribution of a given sequence, but any other function, such as `rbinom()` or custom functions, are valid too.

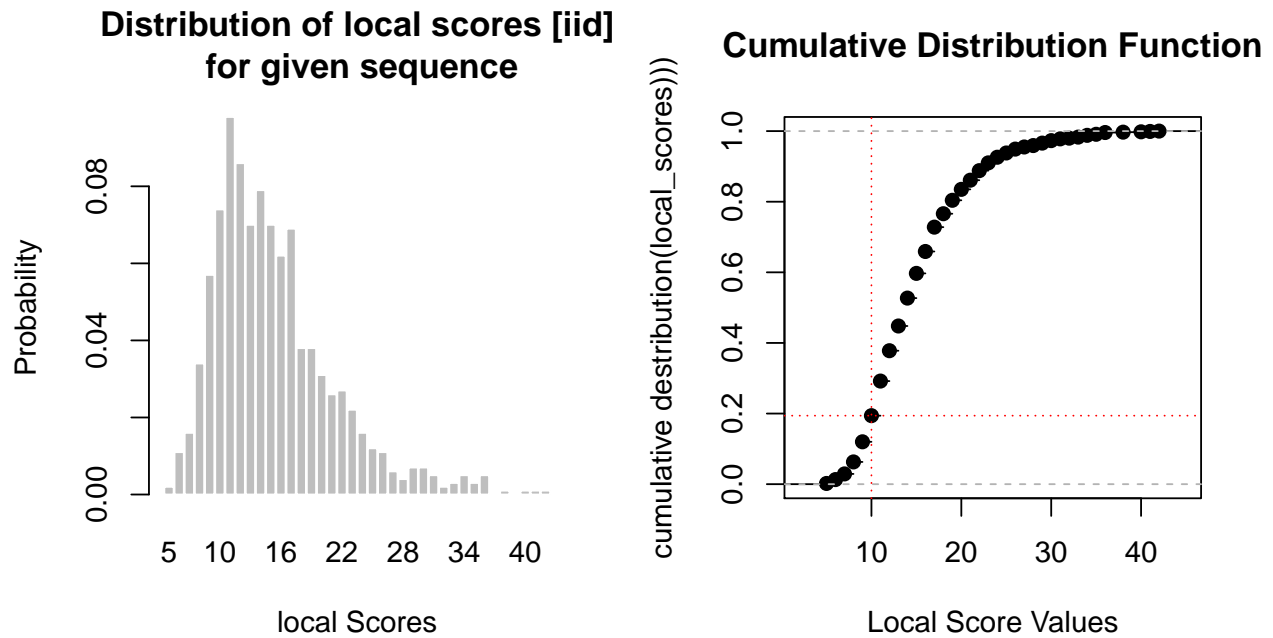
In the following example we search the probability to obtain a local score of 10 in the sequence we created in the previous section and which serves as a blueprint for the score sequences to produce. The return value is the *p*-value of the local score for the given score sequence. A plot of the distribution of all local scores simulated and the cumulative distribution function are displayed. These plots can be hidden by setting the argument “`hist`” to `FALSE`. The number of sequences simulated in our example is 1000, a default value, and can be changed by setting the argument “`numSim`” to an appropriate value.

Note that the cumulative distribution function plot indicates $P(\text{LocalScore} \leq \cdot)$ and so the corresponding *p*-value equals 1 minus the cumulative distribution function value.

```

monteCarlo(local_score = 10, FUN = function(x) {
  return(sample(x = x, size = length(x), replace = TRUE))}, x=sequence)

```



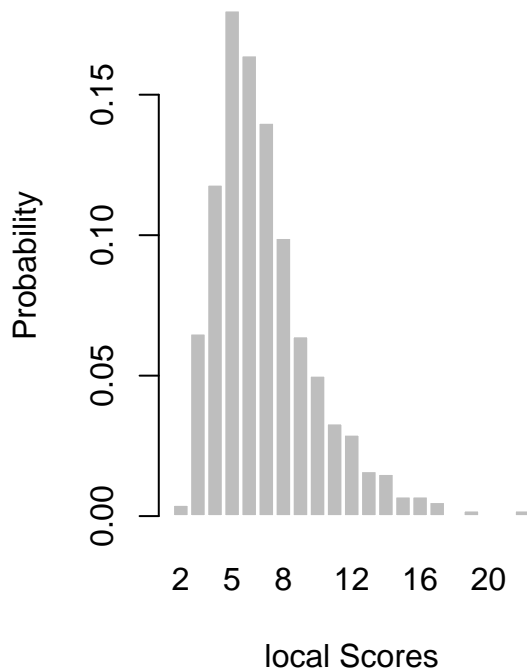
```
#> p-value
#> 0.88
```

The use of this method depends of computing power, number of simulations, implementation of the simulating function and the length of the sequence. For long sequences, you may prefer the next method which combines simulating and approximated methods and called `KarlinMonteCarlo()` (see next section).

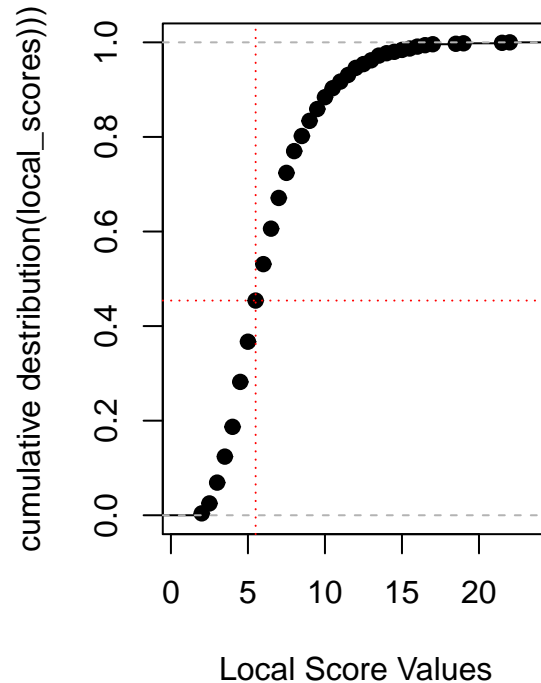
There also exists a special function for real scores.

```
# Example
score_reels=c(-1,-0.5,0,0.5,1)
proba_score_reels=c(0.2,0.3,0.1,0.2,0.2)
sample_from_model <- function(score.sple,proba.sple, length.sple){sample(score.sple,
                                size=length.sple, prob=proba.sple, replace=TRUE)}
monteCarlo_double(5.5,FUN=sample_from_model, plot = TRUE, score.sple=score_reels,proba.sple=proba_score_reels,
                  length.sple=100, numSim = 1000)
```

**Distribution of local scores [iid]
for given sequence**



Cumulative Distribution Function

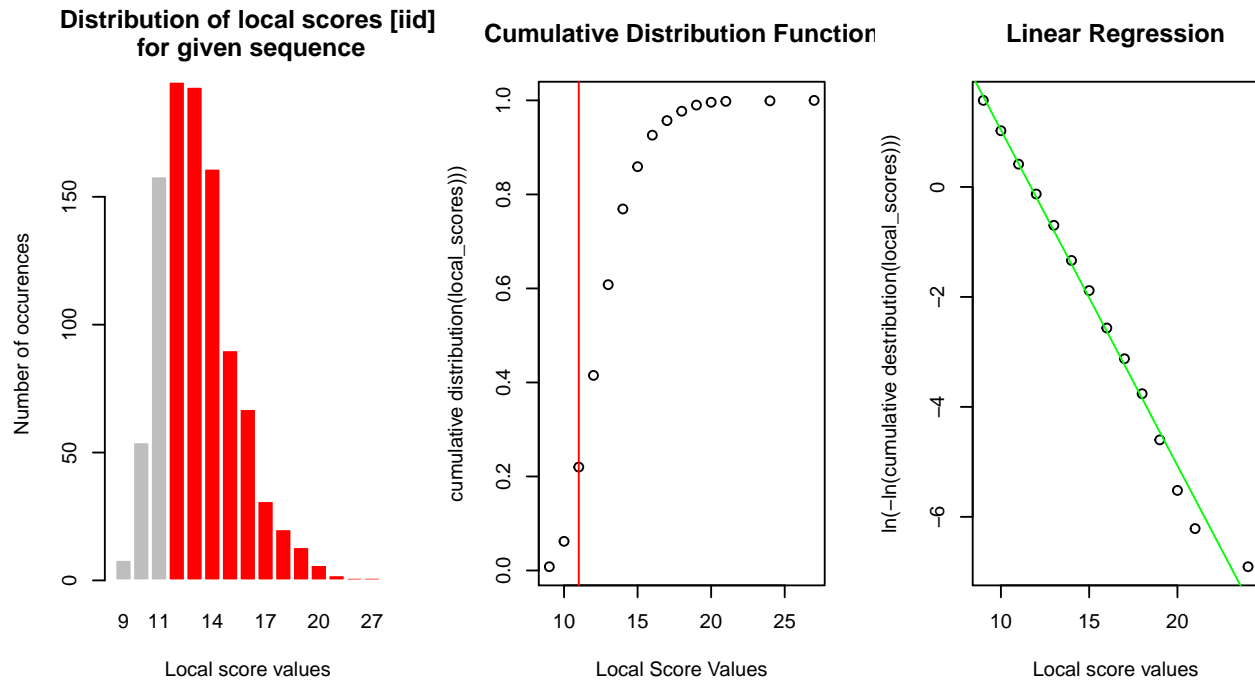


```
#> p-value
#> 0.633
```

A mixed method: functions “karlinMonteCarlo()” and “karlinMonteCarlo_double()”

The function `karlinMonteCarlo()` also uses a function supplied by the user to do simulations. However, it does not deduce directly the p -value from the cumulative distribution function. However, this function is used to estimate the parameters of the Gumbel distribution of *Karlin and al.* approximation. Thus, it is suited for long sequences. Note: `simulated_sequence_length` is the length of the simulated sequences. This value must correspond to the length of sequences yielded by FUN. The value of n is the sequence length for which we want to compute the p -value. If n too large function `MonteCarlo()` could be too much time consuming. Using `karlinMonteCarlo()` with a smaller sequence length `simulated_sequence_length` allows to extract the parameters and then to apply them to the sequence value n .

```
fu<-function(n, size, prob, shift){rbinom(size = size, n = n, prob = prob)+shift}
karlinMonteCarlo(12, FUN = fu, n = 10000, size = 8, prob=0.2, shift = -2,
                 sequence_length = 1000000, simulated_sequence_length = 10000)
```

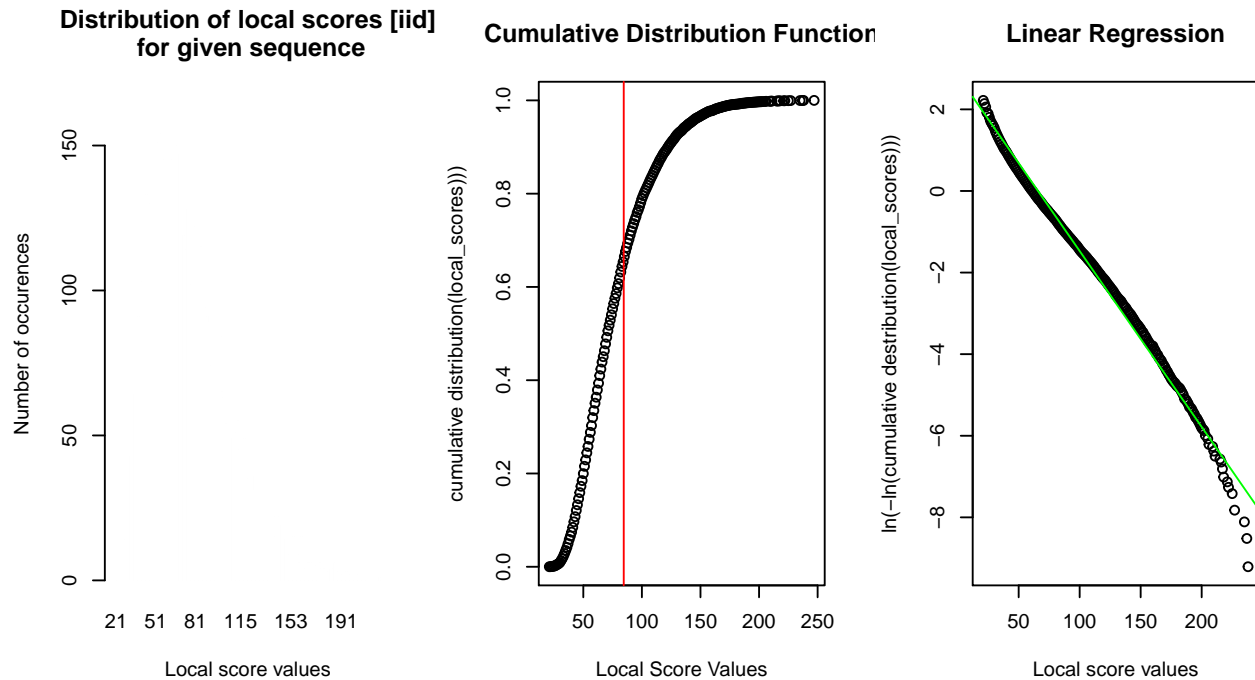


```
#> $`p-value`
#> [1] 1
#>
#> $`K*`
#> [1] 0.1257423
#>
#> $lambda
#> [1] 0.6102375
```

If not specified otherwise, the function produces three graphes, two of them like the function `monte_carlo` and the third a representation of $\ln(-\ln(cf))$, with cf being the cumulated function, showing the linear regression in green color providing the parameters for the Gumbel distribution K^* and λ .

Example for real scores:

```
score_reels=c(-1,-0.5,0,0.5,1)
proba_score_reels=c(0.2,0.3,0.1,0.2,0.2)
fu <- function(score.sple, proba.sple, length.sple){sample(score.sple,
  size=length.sple, prob=proba.sple, replace=TRUE)}
karlinMonteCarlo(85.5, FUN = fu, score.sple=score_reels,proba.sple=proba_score_reels,
  length.sple=10000, numSim = 10000, sequence_length = 100000,
  simulated_sequence_length = 10000)
```

```
#> $`p-value`
#> [1] 0.9862588
#>
#> $`K*`
#> [1] 0.001709635
#>
#> $lambda
#> [1] 0.04310853
```

Exact method for integer scores: function “daudin()”

The exact method calculates the p -value exploiting the fact of the Lindley process associated to the score sequence is a Markov process. Therefore, an exact p -value can be retrieved. The complexity of matrix multiplication involved being $> O(n^2)$, the method is unsuited for sequences of either great length, $n \geq 10^4$ for example could take too much time for computation, or dispersed scores. Note that the exact method requires integer scores.

```
daudin(localScore = 15, sequence_length = 500, score_probabilities =
  c(0.2, 0.3, 0.3, 0.0, 0.1, 0.1), sequence_min = -3, sequence_max = 2)
#> [1] 0.0004119255
```

This function is based on: *S. Mercier and J.J. Daudin 2001: “Exact distribution for the local score of one i.i.d. random sequence”*

How to use the exact method for real scores

The exact method requires integer score to be used. For real scores, a homothetic transformation can be considered to be able to use the exact method. This transformation is theoretically validated to still provide an exact probability. The only existing drawback of the homothetic transformation is that the computation time is increasing. We can check in the following example that the p -value computation using the exact method

with an homothetic transformation corresponds, approximately, to the real local score p -value computed with the Monte Carlo method.

```
score_reels=c(-1,-0.5,0,0.5,1)
proba_score_reels=c(0.2,0.3,0.1,0.2,0.2)
sample_from_model <- function(score.sple,proba.sple, length.sple){sample(score.sple,
                                size=length.sple, prob=proba.sple, replace=TRUE)}
seq.essai=sample_from_model(score.sple=score_reels,proba.sple=proba_score_reels, length.sple=100)
localScoreC_double(seq.essai)
#> $localScore
#> value begin end
#> 10.5 48.0 94.0
#>
#> $suboptimalSegmentScores
#> value begin end
#> [1,] 0.5 3 3
#> [2,] 1.0 8 8
#> [3,] 3.5 12 16
#> [4,] 0.5 46 46
#> [5,] 10.5 48 94
#>
#> $RecordTime
#> [1] 3 8 12 46 48
C=10 # homothetic coefficient
localScoreC(C*seq.essai)
#> $localScore
#> value begin end
#> 105 48 94
#>
#> $suboptimalSegmentScores
#> value begin end
#> [1,] 5 3 3
#> [2,] 10 8 8
#> [3,] 35 12 16
#> [4,] 5 46 46
#> [5,] 105 48 94
#>
#> $RecordTime
#> [1] 3 8 12 46 48
```

We can check that the local score of the sequence which has been homothetically transformed is multiplied by the same coefficient. We are going to compute the p -value. For this we need to create the integer score vector and its corresponding distribution from the ones dedicated to real scores:

```
RealScores2IntegerScores(score_reels,proba_score_reels, coef=C)
#> $ExtendedIntegerScore
#> [1] -10 -9 -8 -7 -6 -5 -4 -3 -2 -1 0 1 2 3 4 5 6 7 8
#> [20] 9 10
#>
#> $ProbExtendedIntegerScore
#> -10 -9 -8 -7 -6 -5 -4 -3 -2 -1 0 1 2 3 4 5 6 7 8 9
#> 0.2 0.0 0.0 0.0 0.0 0.3 0.0 0.0 0.0 0.0 0.1 0.0 0.0 0.0 0.0 0.2 0.0 0.0 0.0 0.0
#> 10
#> 0.2
M.s.r=RealScores2IntegerScores(score_reels,proba_score_reels, coef=C)$ExtendedIntegerScore
```

```

M.s.prob=RealScores2IntegerScores(score_reels,proba_score_reels, coef=C)$ProbExtendedIntegerScore
M.SL=localScoreC(C*seq.essai)$localScore[1]
M.SL
#> value
#> 105
pval.E=daudin(localScore = M.SL,sequence_length = 100, score_probabilities=M.s.prob,
              sequence_min = -10, sequence_max = 10)
pval.E
#> [1] 0.1096559

```

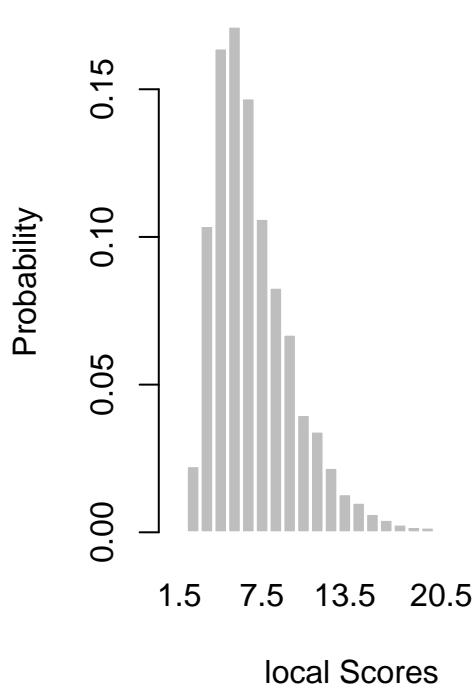
Let us compare the result of the exact method with homothetical transformation with MonteCarlo method for the initial sequence and real scores

```

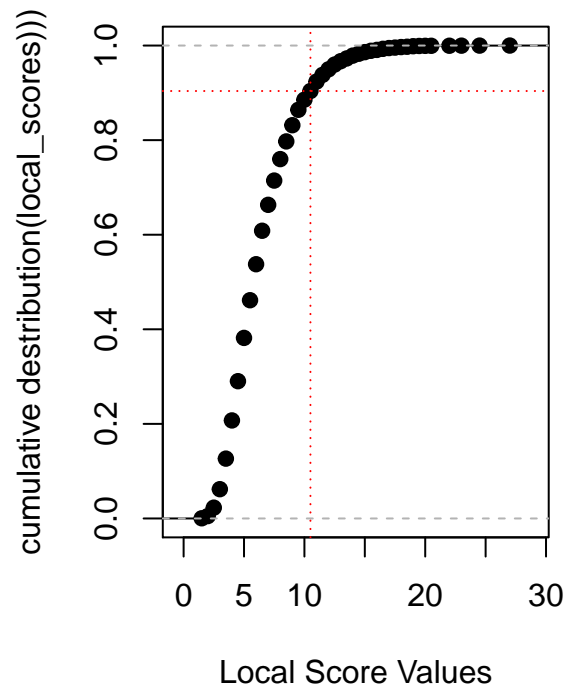
SL.real=localScoreC_double(seq.essai)$localScore[1]
SL.real
#> value
#> 10.5
pval.MC=monteCarlo_double(local_score = SL.real, FUN = sample_from_model,
                          score.sple=score_reels,proba.sple=proba_score_reels, length.sple=100, plot=TRUE,numSim = 10000)

```

**Distribution of local scores [iid]
for given sequence**



Cumulative Distribution Function



```

pval.MC
#> p-value
#> 0.1144

```

We can check that the two probabilities are very similar: 0.1096559 for the exact method and 0.1144 for Monte Carlo's one. The difference comes from the fact that the Monte Carlo method produces an approximation.

Approximate method of Karlin *et al.*: function “karlin()”

The method of Karlin uses the local score’s distinctive cumulative distribution following a law of Gumbel to approximate the p -value. It is suited for large and very large sequences as the approximation is asymptotic with the sequence length and so more accurate for large sequences ; and secondly very large sequence case can be too much time and space consuming for the exact method whereas Karlin *et al.* method does not depend to the sequence length for the computational criteria. The average score must be non positive.

```
score.v=-2:1
score.p=c(0.3, 0.2, 0.2, 0.3)
mean(score.v*score.p)
#> [1] -0.125
karlin(localScore = 14, sequence_length = 100000, sequence_min = -2, sequence_max = 1,
       score_probabilities = c(0.3, 0.2, 0.2, 0.3))
#> [1] 0.3757993
karlin(localScore = 14, sequence_length = 1000, sequence_min = -2, sequence_max = 1,
       score_probabilities = c(0.3, 0.2, 0.2, 0.3))
#> [1] 0.004963226
```

We verify here that the same local score value 14 is more usual for a longer sequence.

```
# With missing score values
karlin(localScore = 14, sequence_length = 1000, sequence_min = -3, sequence_max = 1,
       score_probabilities = c(0.3, 0.2, 0.0, 0.2, 0.3))
#> [1] 0.0006230756
```

This function is based on: Karlin *et al.* 1990: “Methods for assessing the statistical significance of molecular sequence features by using general scoring schemes”

The function `Karlin()` is dedicated for integer scores. For real ones the homothetic solution presented for the exact method can also be applied.

An improved approximate method: function “mcc()”

The function `mcc()` uses an improved version of the Karlin’s method to calculate the p -value. It is suited for sequences of length upper or equal to several hundreds. Let us compare the three methods on the same case.

```
mcc(localScore = 14, sequence_length = 1000, sequence_min = -3, sequence_max = 2,
    score_probabilities = c(0.2, 0.3, 0.3, 0.0, 0.1, 0.1))
#> [1] 0.002423675

daudin(localScore = 14, sequence_length = 1000, score_probabilities =
    c(0.2, 0.3, 0.3, 0.0, 0.1, 0.1), sequence_min = -3, sequence_max = 2)
#> [1] 0.001988438
karlin(localScore = 14, sequence_length = 1000, sequence_min = -3, sequence_max = 2,
    score_probabilities = c(0.2, 0.3, 0.3, 0.0, 0.1, 0.1))
#> [1] 0.0008783872
```

We can observe than the improved approximation method with `mcc()` gives a p -value equal to 0.0024237 which is more accurate than the one of `karlin()` function equal to 8.783872×10^{-4} compared to the exact method which computation equal to 0.0019884.

This function is based on the work of S. Mercier, D. Cellier and D. Charlot 2003 “An improved approximation for assessing the statistical significance of molecular sequence features”*.

An automatic method: function “automatic_analysis()”

This function is meant as a support for the unexperienced user. Since the use of methods for p -value requires some understanding on how these methods work, this function automatically selects an adequate methods based on the sequence given.

There are different use-case scenarios for this function. One can just put a sequence and the model (Markov chains or identically and independantly distributed).

```
automatic_analysis(sequences = list("x1" = c(1,-2,2,3,-2, 3, -3, -3, -2)), model = "iid")
#>
|
|
|
|=====| 100%
#> $x1
#> $x1$p-value`
#> [1] 0.4750898
#>
#> $x1$method applied`
#> [1] "Exact Method Daudin et al"
#>
#> $x1$localScore
#> $x1$localScore$localScore
#> value begin end
#>      6      3      6
#>
#> $x1$localScore$suboptimalSegmentScores
#>      value begin end
#> [1,]      1      1      1
#> [2,]      6      3      6
#>
#> $x1$localScore$RecordTime
#> [1] 1 3
```

In the upper example, the sequence is short, so the exact method is adapted. Here is another example. As the sequence is much more longer, the asymptotic approximation of Karlin *et al.* can be used. This is possible because the average score is negative. If not the MonteCarlo method could have been preferred by the function.

```
score=c(-2,-1,0,1,2)
proba_score=c(0.2,0.3,0.1,0.2,0.2)
sum(score*proba_score)
#> [1] -0.1
sample_from_model <- function(score.sple,proba.sple, length.sple){sample(score.sple,
                                size=length.sple, prob=proba.sple, replace=TRUE)}
seq.essai=sample_from_model(score.sple=score,proba.sple=proba_score, length.sple=5000)
MyAnalysis=automatic_analysis(sequences = list("x1" = seq.essai),
                              distribution=proba_score,score_extremes=c(-2,2), model = "iid")$x1
#>
|
|
|
|=====| 100%
MyAnalysis$p-value`
```

```
#> NULL
MyAnalysis$'method applied'
#> NULL
MyAnalysis$localScore$localScore
#> value begin end
#> 42 2448 2715
```

For real score, achieved the homothetic transformation before. If not, the results could not be correct.

```
score_reels=c(-1,-0.5,0,0.5,1)
proba_score_reels=c(0.2,0.3,0.1,0.2,0.2)
sample_from_model <- function(score.sple,proba.sple, length.sple){sample(score.sple,
                                size=length.sple, prob=proba.sple, replace=TRUE)}
seq.essai=sample_from_model(score.sple=score_reels,proba.sple=proba_score_reels, length.sple=1000)

# Homothetic
RealScores2IntegerScores(score_reels,proba_score_reels, coef=C)
#> $ExtendedIntegerScore
#> [1] -10 -9 -8 -7 -6 -5 -4 -3 -2 -1 0 1 2 3 4 5 6 7 8
#> [20] 9 10
#>
#> $ProbExtendedIntegerScore
#> -10 -9 -8 -7 -6 -5 -4 -3 -2 -1 0 1 2 3 4 5 6 7 8 9
#> 0.2 0.0 0.0 0.0 0.0 0.3 0.0 0.0 0.0 0.0 0.1 0.0 0.0 0.0 0.0 0.2 0.0 0.0 0.0 0.0
#> 10
#> 0.2
M.s.r=RealScores2IntegerScores(score_reels,proba_score_reels, coef=C)$ExtendedIntegerScore
M.s.prob=RealScores2IntegerScores(score_reels,proba_score_reels, coef=C)$ProbExtendedIntegerScore
# The analysis
MyAnalysis=automatic_analysis(sequences = list("x1" = 10*seq.essai), model = "iid",
                              distribution=M.s.prob,score_extremes=range(M.s.r))

#>
|
|
|
|=====| 100%
MyAnalysis$x1$'p-value'
#> [1] 0.9545305
MyAnalysis$x1$'method applied'
#> [1] "Exact Method Daudin et al"

# Without the homothetic, the function gives a wrong result
MyAnalysis2=automatic_analysis(sequences = list("x1" = seq.essai), model = "iid")
#>
|
|
|
|=====| 100%
MyAnalysis2$x1$'p-value'
#> [1] 4.619341e-09
MyAnalysis2$x1$'method applied'
#> [1] "Exact Method Daudin et al"
```

Whenever there is no distribution given, these is learnt from the sequence(s) given. The sequence(s) must

be passed as **named list**. If there are more than one sequence, all sequences are treated. A progress bar informs the user about the progress made (helpful for computational intense sequences and/or great numbers of sequences). If the sequence(s) passed are not score sequences but sequence(s) of letters, a file picker dialog pops up. Here, the user can choose a file containing a mapping from letters to scores. The format is csv (view section “File Formats” for details) and allows also to provide a distribution that is loaded concurrently to the score. One can also chose not to provide sequences at all. In this case, the first dialog that pops up is for selection of a FASTA file (view section “File Formats” for details). Either way, the user has little influence on the method applied to find the p -value. One can change the argument “method_limit” to modify the threshold for the use of exact and approximating methods or supply a function for simulation methods which will result in the use of a simulation method. In this case, make sure to provide a suitable “simulated_sequence_length” argument, as for long sequences, the method “monteCarloKarlin” will be used.

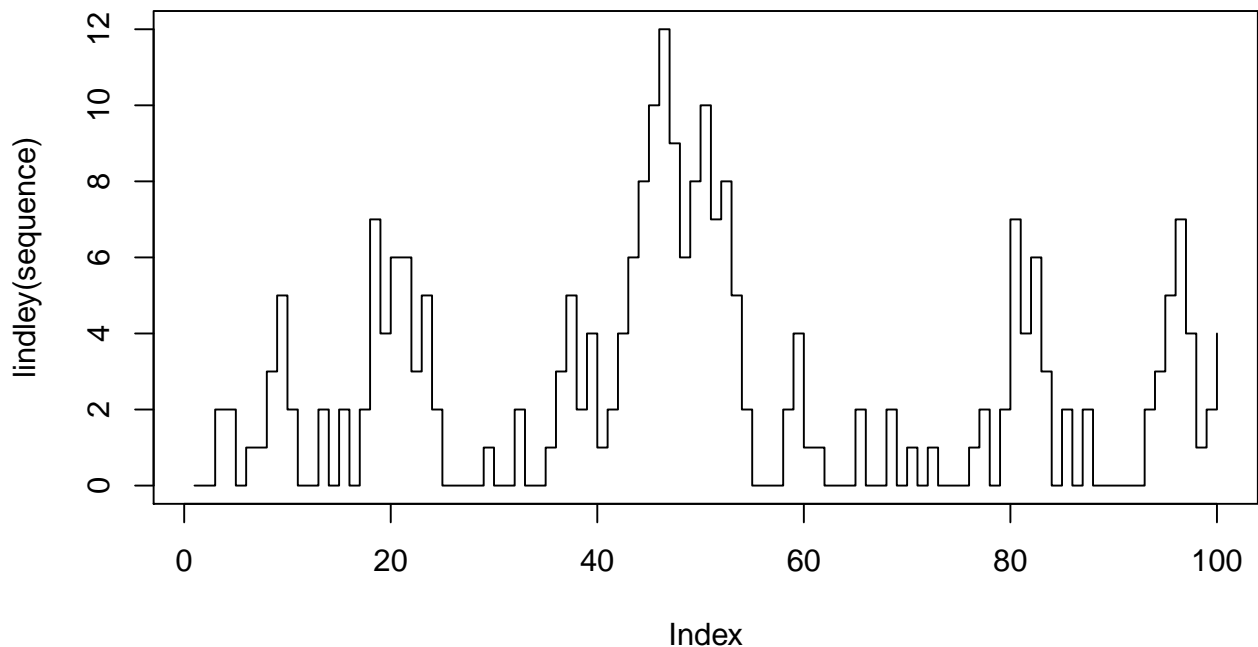
Other Functions

The package provides a number of auxiliary functions for loading or creating sequences for different models.

Lindley Process: to visualize optimal and suboptimal segments

This function takes a score sequence and shows the Lindley process associated, illustrating the operation of the local score algorithm. Plotting the Lindley process provides a view of the potential optimal segments and their comparison between them.

```
set.seed(1)
sequence <- sample(c(-3,2,0,1,5), 100, replace = TRUE, prob = c(0.5, 0.3, 0.05, 0.1, 0.05))
lindley(sequence)
#> [1] 0 0 2 2 0 1 1 3 5 2 0 0 2 0 2 0 2 7 4 6 6 3 5 2 0
#> [26] 0 0 0 1 0 0 2 0 0 1 3 5 2 4 1 2 4 6 8 10 12 9 6 8 10
#> [51] 7 8 5 2 0 0 0 2 4 1 1 0 0 0 2 0 0 2 0 1 0 1 0 0 0
#> [76] 1 2 0 2 7 4 6 3 0 2 0 2 0 0 0 0 0 2 3 5 7 4 1 2 4
plot(lindley(sequence), type = "s")
```



Score Loading Function

`loadScoreFromFile()` reads a csv-file returning a named list where names correspond to the first file column and values to the second file column. If a third column is available within the file, it will be read and can be accessed by name, too. An example is given in the first case study. See for example the provided files `Kyte1982.txt` or `score_file_example.csv`. The function is reading a header line by default, so be careful that you have one otherway the first score will be missed.

Empirical distribution: function “`scoreSequences2probabilityVector()`”

`scoreSequences2probabilityVector()` takes in a list of score sequences and returns the resulting empirical distribution from the minimal to the maximal score value of all sequences. Thus,

```
seq1 = sample(7:8, size = 10, replace = TRUE)
seq2 = sample(2:3, size = 15, replace = TRUE)
l = list(seq1, seq2)
r = scoreSequences2probabilityVector(l)
r
#>      2      3      4      5      6      7      8
#> 0.36 0.24 0.00 0.00 0.00 0.20 0.20
length(r)
#> [1] 7
```

returns a vector of length 7, even if there are only 4 distinct unique values present in the list. Very useful for the use in any non-simulating method of p -value.

Case study

Medium sequence

Sequence extracted from <https://www.uniprot.org/uniprot/P49755/>

```
data(MidSeq)
MySeq=MidSeq
data(dico)
SeqScore <- CharSequence2ScoreSequence(MySeq,dico)
n <- length(SeqScore)
n
#> [1] 219
```

Local score computation and parameter model settings

```
LS <- localScoreC(SeqScore)$localScore[1]
LS
#> value
#>      52
```

Parameter model settings


```

prob = scoreSequences2probabilityVector(list(SeqScore))
prob
#>      -5      -4      -3      -2      -1      0      1
#> 0.06392694 0.25570776 0.02283105 0.03652968 0.15068493 0.06849315 0.00000000
#>      2      3      4      5
#> 0.08675799 0.07762557 0.17808219 0.05936073

```

Exact method

```

time.daudin <- system.time(
res.daudin <- daudin(localScore = LS, sequence_length = n,
  score_probabilities = prob,
  sequence_min = min(SeqScore),
  sequence_max = max(SeqScore)))
res.daudin
#> [1] 0.2654051

```

Approximated method

The call of the function `karlin()` is similar to the one of `daudin()`.

```

time.karlin <- system.time(
res.karlin <- karlin(localScore = LS, sequence_length = n,
  score_probabilities = prob,
  sequence_min = min(SeqScore),
  sequence_max = max(SeqScore)))
res.karlin
#> [1] 0.1964297

```

The two p -values are different because the sequence length n is equal 219. It is not enough large to have a good approximation with the approximated method.

Improved approximation

The call of the function ‘`mcc()`’ is still the same.

```

time.mcc <- system.time(
  res.mcc <- mcc(localScore = LS, sequence_length = n,
    score_probabilities = prob,
    sequence_min = min(SeqScore),
    sequence_max = max(SeqScore)))
res.mcc
#> [1] 0.7773828

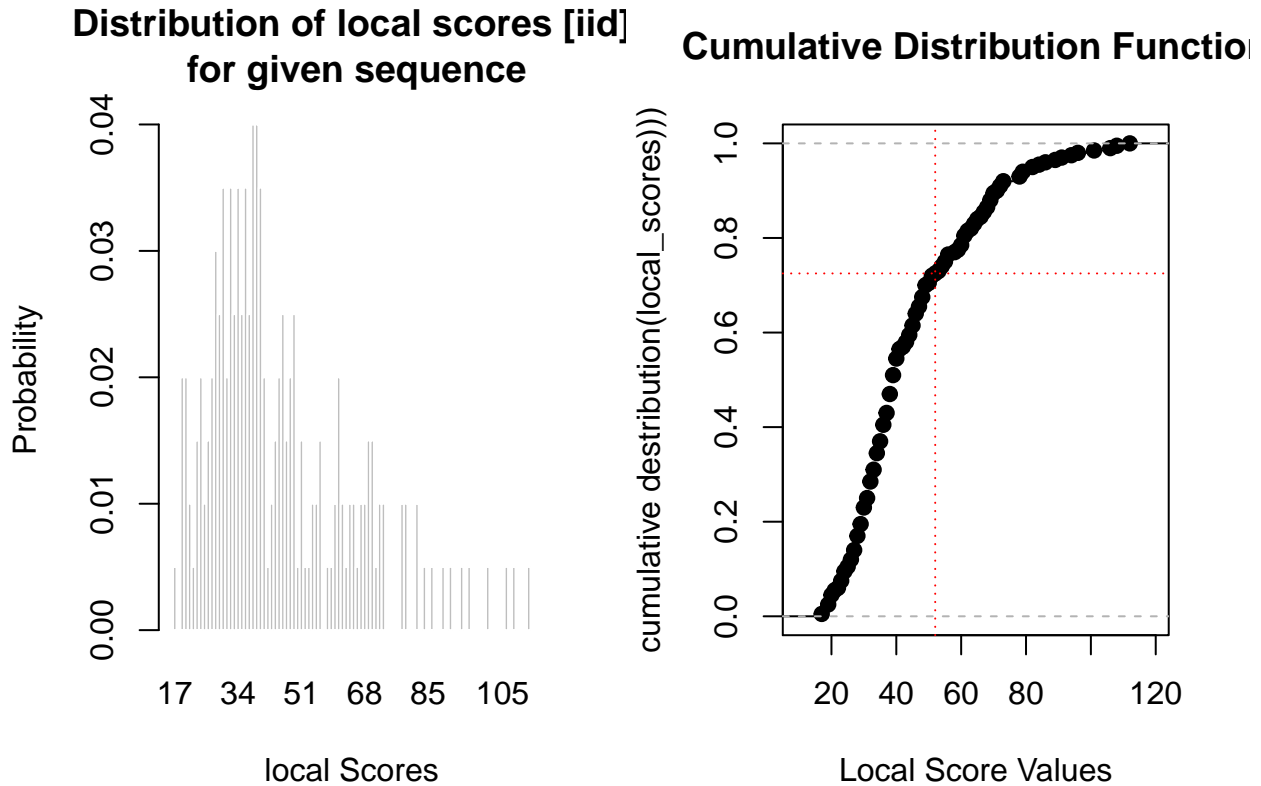
```

We can verify here that this approximation is more accurate than the one of Karlin *et al.* for sequences of length of several hundred components.

Monte Carlo

Let us do a very quick empirical computation with only 200 repetitions.

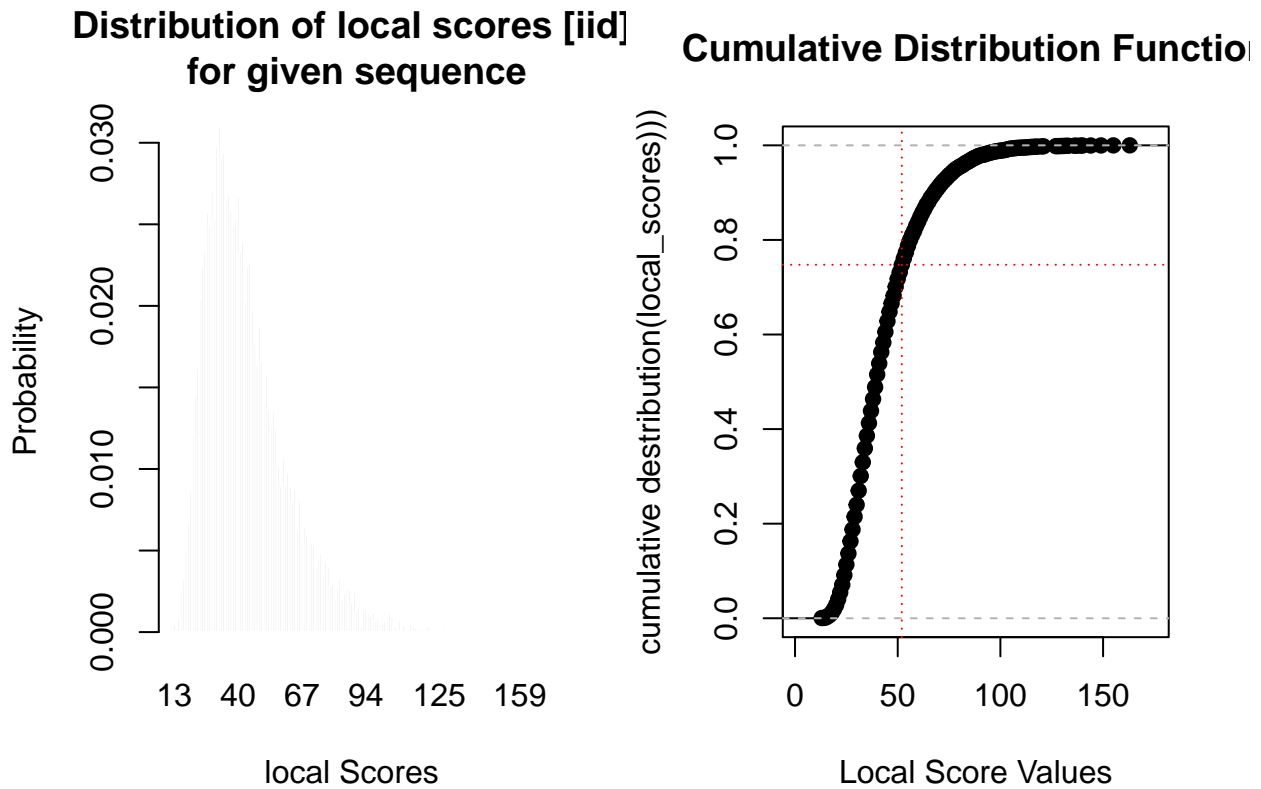
```
time.MonteCarlo1 <- system.time(
res.MonteCarlo1 <- monteCarlo(local_score = LS,
  FUN = function(x) {return(sample(x = x,size = length(x),
    replace = TRUE))},
  x=SeqScore, numSim = 200))
```



```
res.MonteCarlo1
#> p-value
#> 0.28
```

The p -value estimation is 0.28 which is around the exact value 0.2654051. Let us increase the number of repetition to be more accurate.

```
time.MonteCarlo2 <- system.time(
res.MonteCarlo2 <- monteCarlo(local_score = LS,
  FUN = function(x) {return(sample(x = x,size = length(x),
    replace = TRUE))},
  x=SeqScore, numSim = 10000))
```



```
res.MonteCarlo2
#> p-value
#> 0.2683
```

Result and time computation comparison

P-value

```
res.pval <- c(Daudin = res.daudin, Karlin = res.karlin, MCC = res.mcc,
  MonteCarlo1=res.MonteCarlo1, MonteCarlo2=res.MonteCarlo2)
names(res.pval) = c("Exact", "Approximation", "Improved appx", "MonteCarlo1", "MonteCarlo2")
res.pval
#>      Exact Approximation Improved appx MonteCarlo1 MonteCarlo2
#> 0.2654051 0.1964297 0.7773828 0.2800000 0.2683000
```

Computation time

```
rbind(time.daudin, time.karlin, time.mcc, time.MonteCarlo1, time.MonteCarlo2)
#>      user.self sys.self elapsed user.child sys.child
#> time.daudin      0.001      0 0.000      0      0
#> time.karlin      0.000      0 0.000      0      0
#> time.mcc         0.000      0 0.000      0      0
#> time.MonteCarlo1 0.018      0 0.018      0      0
#> time.MonteCarlo2 0.179      0 0.179      0      0
```

Time computation for MonteCarlo method depends to the number of repetition. For medium sequences exact method is around 30 time more time cumsumming than the mcc method. Karlin's method is the fastest one but can be not accurate if the sequence are too short (here $n = 219$ a couple of hundred is note enough, a thousand may be preferred to be sure of the accuracy).

Short sequence

```
data(ShortSeq)
MySeq.Short =ShortSeq
SeqScore.Short <- CharSequence2ScoreSequence(MySeq.Short,dico)
n.short <- length(SeqScore.Short)
n.short
#> [1] 31
```

Sequence of length 31. For short sequences, it is easier and usual to obtain a positive expectation for the score. So the functions based on approximated methods can't be used and an error message is given.

```
SeqScore.S <- SeqScore.Short
LS.S <- localScoreC(SeqScore.S)$localScore[1]
prob.S = scoreSequences2probabilityVector(list(SeqScore.S))
```

```
LS.S
#> value
#> 52
prob.S
#>      -5      -4      -3      -2      -1      0      1
#> 0.03225806 0.06451613 0.00000000 0.00000000 0.22580645 0.06451613 0.00000000
#>      2      3      4      5
#> 0.09677419 0.09677419 0.25806452 0.16129032
```

```
time.daudin <- system.time(
res.daudin. <- daudin(localScore = LS.S, sequence_length = n.short,
  score_probabilities = prob.S,
  sequence_min = min(SeqScore.S),
  sequence_max = max(SeqScore.S)))

time.karlin <- system.time(
res.karlin <- try(karlin(localScore = LS.S, sequence_length = n.short,
  score_probabilities = prob.S,
  sequence_min = min(SeqScore.S),
  sequence_max = max(SeqScore.S))))
#> Error in karlin(localScore = LS.S, sequence_length = n.short, score_probabilities = prob.S, :
#> [Invalid Input] Score expectation must be strictly negative.
```

```
time.mcc <- system.time(
  res.mcc <- try(mcc(localScore = LS.S, sequence_length = n.short,
    score_probabilities = prob.S,
    sequence_min = min(SeqScore.S),
    sequence_max = max(SeqScore.S))))
#> Error in mcc(localScore = LS.S, sequence_length = n.short, score_probabilities = prob.S, :
#> [Invalid Input] Score expectation must be strictly negative.
```

```
time.karlinMonteCarlo <- system.time(
res.karlinMonteCarlo <-
  karlinMonteCarlo(local_score = LS.S, plot=FALSE,
    sequence_length = n.short,
    simulated_sequence_length = 1000,
    FUN = sample, x=min(SeqScore.S):max(SeqScore.S),
    size = 1000, prob=prob.S, replace=TRUE,
```

```

numSim = 100000))

time.MonteCarlo <- system.time(
res.MonteCarlo <- monteCarlo(local_score = LS.S, plot=FALSE,
  FUN = function(x) {return(sample(x = x, size = length(x),
    replace = TRUE))},
  x=SeqScore.S, numSim = 10000))

```

Results

```

res.pval <- c(Daudin = res.daudin, MonteCarlo=res.MonteCarlo)
names(res.pval) = c("Daudin", "MonteCarlo")
res.pval
#>      Daudin MonteCarlo
#> 0.2654051 0.5895000
rbind(time.daudin, time.MonteCarlo)
#>      user.self sys.self elapsed user.child sys.child
#> time.daudin      0.000      0 0.001      0      0
#> time.MonteCarlo 0.101      0 0.101      0      0

```

Here an example using another probability vector with non positive average score. In this example, the local score is very huge and realized by the whole sequence, the p -value is very low as confirmed by the exact method.

```

set.seed(1)
prob.bis=dnorm(-5:5, mean=-0.5, sd=1)
prob.bis=prob.bis/sum(prob.bis)
names(prob.bis)=-5:5
# Score Expectation
sum((-5:5)*prob.bis)
#> [1] -0.4999994

time.mcc <- system.time(
  res.mcc <- mcc(localScore = LS.S, sequence_length = n.short,
    score_probabilities = prob.bis,
    sequence_min = min(SeqScore.S),
    sequence_max = max(SeqScore.S)))

time.daudin <- system.time(
res.daudin <- daudin(localScore = LS.S, sequence_length = n.short,
  score_probabilities = prob.bis,
  sequence_min = min(SeqScore.S),
  sequence_max = max(SeqScore.S)))

simu=function(n,p){return(sample(x=-5:5,size = n, replace=TRUE, prob=p))}
time.MonteCarlo <- system.time(
res.MonteCarlo <-
monteCarlo(local_score = LS.S, plot=FALSE,
  FUN = simu, n.short, prob.bis, numSim = 100000))

res.pval <- c(MCC=res.mcc,Daudin = res.daudin, MonteCarlo=res.MonteCarlo)
names(res.pval) = c("MCC", "Daudin", "MonteCarlo")
res.pval

```

```

#>          MCC          Daudin  MonteCarlo
#> 0.000000e+00 1.306309e-33 0.000000e+00
rbind(time.mcc,time.daudin, time.MonteCarlo)
#>          user.self sys.self elapsed user.child sys.child
#> time.mcc          0.000          0 0.000          0          0
#> time.daudin        0.000          0 0.000          0          0
#> time.MonteCarlo    1.095          0 1.095          0          0

```

For short sequences, exact method is fast, more precise and must be preferred.

Large sequence

```

data(LongSeq)
MySeq.Long=LongSeq
SeqScore.Long <- CharSequence2ScoreSequence(MySeq.Long,dico)
n.Long <- length(SeqScore.Long)
n.Long
#> [1] 1093
SeqScore.Long <- CharSequence2ScoreSequence(MySeq.Long,dico)
LS.L <- localScoreC(SeqScore.Long)$localScore[1]
LS.L
#> value
#> 65
prob.L = scoreSequences2probabilityVector(list(SeqScore.Long))
prob.L
#>          -5          -4          -3          -2          -1          0          1
#> 0.07410796 0.20311070 0.02012809 0.07502287 0.21225984 0.07776761 0.00000000
#>          2          3          4          5
#> 0.07136322 0.09423605 0.14364135 0.02836231
sum(prob.L*as.numeric(names(prob.L)))
#> [1] -0.4638609

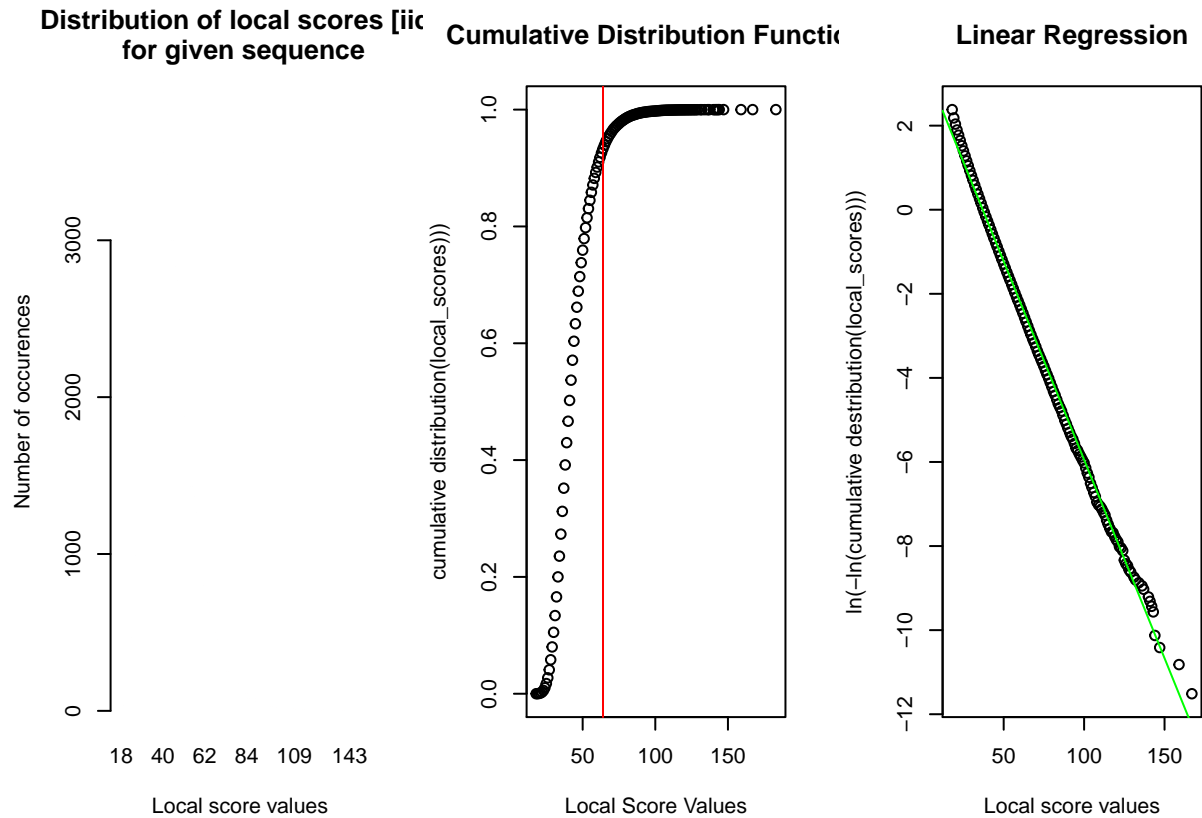
```

Sequence of length 1093 with a local score equal to 65. The average score is non positive so approximated methods can be used.

```

#> [1] 0.07231933
#> p-value
#> 0.07313

```



```
#> $`p-value`
#>      value
#> 0.07467078
#>
#> $`K*`
#> [1] 0.03279319
#>
#> $lambda
#> [1] 0.09438864
```

Results

```
res.pval.L <- c(res.daudin.L, res.mcc.L, res.karlin.L, res.karlinMonteCarlo.L$`p-value`,res.MonteCarlo.L)
names(res.pval.L) = c("Daudin","MCC","Karlin","KarlinMonteCarlo","MonteCarlo")
res.pval.L
#>      Daudin      MCC      Karlin KarlinMonteCarlo
#> 0.07231933 0.22338343 0.07639838 0.07467078
#> MonteCarlo
#> 0.07313000

rbind(time.daudin.L, time.karlin.L, time.mcc.L,time.karlinMonteCarlo.L,
      time.MonteCarlo.L)
#>      user.self sys.self elapsed user.child sys.child
#> time.daudin.L      0.001      0    0.001      0      0
#> time.karlin.L      0.000      0    0.001      0      0
#> time.mcc.L        0.000      0    0.000      0      0
#> time.karlinMonteCarlo.L 4.537      0    4.537      0      0
```

```
#> time.MonteCarlo.L          4.390          0  4.390          0          0
```

Even for large sequences of several thousands, the exact method is still fast enough but it could become too much time consuming for a sequence data set with numerous sequences. The approximated methods must be preferred.

Several sequences

The function `automatic_analysis()` can analysis a named list of sequences. It choose the adequate method for each sequence. Here in the following example, the exact method is used for the short sequence, whereas an asymptotic method is used for the long one.

```
MySeqsList=list(MySeq,MySeq.Short,MySeq.Long)
names(MySeqsList)=c('Q09FU3.fasta','P49755.fasta','Q60519.fasta')
MySeqsScore=lapply(MySeqsList, FUN=CharSequence2ScoreSequence, dico)
AA=automatic_analysis(MySeqsScore, model='iid')
#>
|
|
|=====| 33%
|=====| 67%
|=====| 100%
AA$Q09FU3.fasta
#> $`p-value`
#> [1] 0.08626993
#>
#> $`method applied`
#> [1] "Exact Method Daudin et al"
#>
#> $localScore
#> $localScore$localScore
#> value begin end
#> 52 14 38
#>
#> $localScore$suboptimalSegmentScores
#> value begin end
#> [1,] 5 1 4
#> [2,] 2 9 9
#> [3,] 52 14 38
#> [4,] 17 121 124
#> [5,] 7 138 139
#> [6,] 4 144 144
#> [7,] 4 147 147
#> [8,] 4 149 149
#> [9,] 4 151 151
#> [10,] 4 154 154
#> [11,] 4 157 157
#> [12,] 9 161 162
#> [13,] 2 175 175
#> [14,] 43 186 208
```



```
#>
#> $localScore$RecordTime
#> [1] 1 9 14 121 138 144 147 149 151 154 157 161 175 186
AA$Q09FU3.fasta$`method applied`
#> [1] "Exact Method Daudin et al"
AA$Q60519.fasta$`method applied`
#> [1] "Exact Method Daudin et al"
```

We can observe differences between the p -value of the short sequence obtained in the case study for the only short sequence, and the one obtained with the automatic analysis. Note that the distribution vector of the scores used are different which induces a different p -value.

```
cbind(prob, prob.S, prob.L, '3 sequences'=scoreSequences2probabilityVector(MySeqsScore))
#>      prob      prob.S      prob.L 3 sequences
#> -5 0.06392694 0.03225806 0.07410796 0.07148176
#> -4 0.25570776 0.06451613 0.20311070 0.20848846
#> -3 0.02283105 0.00000000 0.02012809 0.02010424
#> -2 0.03652968 0.00000000 0.07502287 0.06701415
#> -1 0.15068493 0.22580645 0.21225984 0.20253165
#> 0 0.06849315 0.06451613 0.07776761 0.07594937
#> 1 0.00000000 0.00000000 0.00000000 0.00000000
#> 2 0.08675799 0.09677419 0.07136322 0.07446016
#> 3 0.07762557 0.09677419 0.09423605 0.09158600
#> 4 0.17808219 0.25806452 0.14364135 0.15189873
#> 5 0.05936073 0.16129032 0.02836231 0.03648548
```

Using the probability vector of the three sequences to compute the p -value of the local score of the short sequence with the function `daudin()`, we recover an identical p -value than we have obtained with the automatic analysis().

```
daudin.bis=daudin(localScore=LS.S,sequence_length = n.short, score_probabilities = scoreSequences2probabilityVector(MySeqsScore))
daudin.bis
#> [1] 0.0005164345
AA$P49755.fasta$`p-value`
#> [1] 0.0005164345

# automatic_analysis(sequences=list('MySeq.Short'=MySeq.Short), model='iid', distribution=proba.S)
```

A larger example with a SCOP data base

```
library(localScore)
data(dico)
data(MySeqList)
MySeqScoreList=lapply(MySeqList, FUN=CharSequence2ScoreSequence, dico)
head(MySeqScoreList)
#> $P50456
#> [1] 2 -5 -4 4 5 -4 4 4 5 -4 -3 -4 4 0 2 0 4 5 -1 -4 0 -3 4 4 -3
#> [26] 2 0 -1 -1 -1 4 4 -4 5 0 -3 -1 -4 4 -4 -2 -1 0 -4 -5 3 -1 3 0 -4
#> [51] -3 0 3 4 -4 -1 5 2 -1 4 -4 -1 5 4 -4 4 2 -4 4 -5 4 -4 -4 -1 2
#> [76] -1 -1 2 4 -3 0 -4 -2 4 -1 4 -4 -1 4 3 -4 2 2 4 -5 0 -4 4 4 2
#> [101] -4 -4 5 5 -1 0 4 0 2 -3 4 0 -5 5 4 2 5 2 4 -4 4 3 -4 -2 -4
#> [126] -4 5 4 5 0 -1 -2 4 -1 -4 2 2 -4 5 4 3 -2 4 5 -1 -4 -1 5 -5 -4
#> [151] -4 2 4 -2 2 -1 -1 -4 -3 5 -1 4 -4 -1 -1
```

```

#>
#> $P14859
#> [1] -4 -4 -2 -1 -4 4 -4 -4 4 -4 -4 3 2 -4 -1 3 -4 -4 -5 -5 5 -4 4 0 3
#> [26] -1 -4 0 -4 4 0 4 2 2 0 -4 4 -1 0 -4 -4 3 -1 -4 -1 -1 5 -1 -5 3
#> [51] -4 2 4 -4 4 -1 3 -4 -4 2 3 -4 4 -4 -2 4 4 -4 -4 -1 4 -4 -4 2 -4
#>
#> $P14859
#> [1] 5 -4 -1 -4 5 -5 4 2 4 -4 -4 -1 3 4 -4 -4 -4 -4 -2 -1 -1 -4 -4 5 -1
#> [26] 2 5 2 -4 -4 4 -4 2 -4 -4 -4 4 5 -5 4 -1 3 3 -4 -5 -5 -4 -4 -4 -4
#> [51] -5 5
#>
#> $P10037
#> [1] 5 -1 5 2 2 -4 -4 2 4 -4 -5 -3 3 0 -4 -3 -1 -4 -2 -1 -1 -4 -4 5 2
#> [26] -5 2 2 -4 -4 4 -4 4 -4 -4 -4 4 4 -5 4 -1 3 3 -4 -5 -5 -4 -5 -4 -4
#> [51] -5 4
#>
#> $Q13619
#> [1] -4 -1 4 -4 -4 -4 4 -1 -1 -1 -4 -5 4 3 -4 -4 -5 -4 -1 -4 5 -4 2 2 5
#> [26] 4 -5 5 2 -4 2 -5 -4 -1 4 0 -3 -4 4 4 4 -1 -4 4 -1 -4 -4 4 -4 3
#> [51] -2 4 -4 -2 0 -4 4 -4 -4 -5 5 -4 -1 4 5 -4 -5 -4 -1 2 -4 -5 -4 -4 -4
#> [76] -4 -2 -4 -4 -1 -3 -1 4 2
#>
#> $Q13619
#> [1] 4 -1 4 3 -4 -1 4 4 4 4 2 3 -4 -4 0 -4 0 3 -1 3 -4 -4 5 -4 2
#> [26] 2 -1 0 5 -4 -4 -1 -4 4 -5 -5 -1 4 -4 -1 4 2 3 0 -4 2 -5 4 4 5
#> [51] -4 -1 -2 -4 0 -4 -4 4 -4 -4 0 -4 -4 3 5 3 -4
AA=automatic_analysis(sequences=MySeqScoreList, model='iid')
#>
|
|
|
|=
|=
|==
|==
|==
|===
|===
|====
|====
|=====
|

```

```

0%
1%
1%
2%
2%
3%
4%
4%
5%
5%
6%
7%

```

=====	8%
=====	8%
=====	9%
=====	9%
=====	10%
=====	11%
=====	11%
=====	12%
=====	12%
=====	13%
=====	14%
=====	15%
=====	15%
=====	16%
=====	16%
=====	17%
=====	18%
=====	18%
=====	19%
=====	19%
=====	20%
=====	21%
=====	21%
=====	22%
=====	22%
=====	23%
=====	24%

=====	24%
=====	25%
=====	25%
=====	26%
=====	27%
=====	28%
=====	28%
=====	29%
=====	29%
=====	30%
=====	31%
=====	31%
=====	32%
=====	32%
=====	33%
=====	34%
=====	35%
=====	35%
=====	36%
=====	36%
=====	37%
=====	38%
=====	38%
=====	39%
=====	39%
=====	40%

=====	41%
=====	41%
=====	42%
=====	42%
=====	43%
=====	44%
=====	44%
=====	45%
=====	45%
=====	46%
=====	47%
=====	48%
=====	48%
=====	49%
=====	49%
=====	50%
=====	51%
=====	51%
=====	52%
=====	52%
=====	53%
=====	54%
=====	55%
=====	55%
=====	56%
=====	56%
=====	57%

	58%
	58%
	59%
	59%
	60%
	61%
	61%
	62%
	62%
	63%
	64%
	64%
	65%
	65%
	66%
	67%
	68%
	68%
	69%
	69%
	70%
	71%
	71%
	72%
	72%
	73%

=====	74%
=====	75%
=====	75%
=====	76%
=====	76%
=====	77%
=====	78%
=====	78%
=====	79%
=====	79%
=====	80%
=====	81%
=====	81%
=====	82%
=====	82%
=====	83%
=====	84%
=====	84%
=====	85%
=====	85%
=====	86%
=====	87%
=====	88%
=====	88%
=====	89%
=====	89%
=====	90%

```

|=====| 91%
|=====| 91%
|=====| 92%
|=====| 92%
|=====| 93%
|=====| 94%
|=====| 95%
|=====| 95%
|=====| 96%
|=====| 96%
|=====| 97%
|=====| 98%
|=====| 98%
|=====| 99%
|=====| 99%
|=====| 100%
AA[[1]]
#> $`p-value`
#> [1] 0.06389172
#>
#> $`method applied`
#> [1] "Exact Method Daudin et al"
#>
#> $localScore
#> $localScore$localScore
#> value begin end
#> 67 4 144
#>
#> $localScore$suboptimalSegmentScores
#> value begin end
#> [1,] 2 1 1
#> [2,] 67 4 144
#>
#> $localScore$RecordTime
#> [1] 1 4
# the p-value of the first 10 sequences
apply(AA, function(x){x$`p-value`})[1:10]

```



```

#>      P50456      P14859      P10037      Q13619      P22262      P20823      P07014
#> 0.06389172 0.97260747 0.87470986 0.89625648 0.45058860 0.96651306 0.74891930
#>      Q9X399      Q0SB06      Q9I641
#> 0.68145292 0.99374006 0.51203351
# the 20th smallest p-values
sort(sapply(AA, function(x){x$`p-value`}))[1:20]
#>      Q5SMG8      POA334      Q2W6R1      027564      P12282      P50456
#> 9.485100e-07 3.442818e-04 4.406208e-04 4.548065e-04 6.167591e-02 6.389172e-02
#>      Q58194      Q8AA93      P05523      P28793      P55038      Q9WZ12
#> 7.290625e-02 9.038738e-02 9.064666e-02 9.414140e-02 9.498902e-02 1.017792e-01
#>      Q9KQJ1      Q5ZSV0      POA544      P77072      POA9G8      P08709
#> 1.304836e-01 1.341404e-01 1.356180e-01 1.410952e-01 1.481435e-01 1.490303e-01
#>      P00390      Q13564
#> 1.509551e-01 1.574591e-01
which(sapply(AA, function(x){x$`p-value`})<0.05)
#> Q2W6R1 027564 POA334 Q5SMG8
#>      14      90      150      192
table(sapply(AA, function(x){x$`method`}))
#>
#> Exact Method Daudin et al
#>      206
# The maximum sequence length equals 404 so it here normal that the exact method is used for all the 60
scoreSequences2probabilityVector(MySeqScoreList)
#>      -5      -4      -3      -2      -1      0      1
#> 0.05537938 0.26383764 0.02153482 0.04105166 0.14754382 0.07195572 0.00000000
#>      2      3      4      5
#> 0.10487777 0.05241006 0.17481550 0.06659363

```

File Formats

This package allows input in file form. For the package to work, please respect the following conventions.

Sequence Files

The package accepts files in FASTA format: Every sequence is preceeded by a title (marked by a ">") and a line break. One sequence takes one line, followed by a line break and a line only containing a tab.

```

>HUMAN_NM_018998_2
TGAGTAGGGCTGGCAGAGCTGGGGCCTCATGGCTGTGTAGTAGCAGGCCCCGCCCCGCGACCTGGCCAGGCGATCACTACAGCCGCCCTGCCGAACAG

>Mouse_NM_013908_3
CCCCATGAGGACCCAGAACCCTCAATGGAGAAGAGTCAGGATTGTCTGTGCTGCCAGAGTGAAGTGGCCTGGTAATTACCCTGCAGCCTTTCTGGAACAG

>HUMAN_NM_018998_3
GTGAGCACGGGCGGCGGGTTGACCCTGCCCCGCCCCACGCCGACAGCCTGTCCAGCCCCGGCCTCCCCACAG

>Mouse_NM_013908_4
GTAAGTGTGGGCATTGGGTGGGCTACCTGTCCCATTGTGCCCTGCCAGCAGTCTGCCAGCTGTGGCCTTCCCCCAG

>HUMAN_NM_018998_5
GGTGCTCACAGCCCAGAGACACCACTGAGGTAGGAAGCTGCCCTGGAGTGATGTCCTTGGGGCATTGGACAGGGACCCTCACCGTAGCCCTCCCTGCAG

```

Score Files

A score file is a csv file that contains a header line and each line contains a letter and its score. Optionnally one can also provide a probability for each score. Example:

```
Letters,Scores,Probabilities
L,-2,0.04
M,-1,0.04
N,0,0.04
```

Transition Matrix Files

A csv file only containing the values of the matrix. Example:

```
0.2,0.3,0.5
0.3,0.4,0.3
0.2,0.4,0.4
```

A word for Markovian model

At the present time, a markovian dependance of the components of the sequence can be taken into account using the `automatic.analysis()` or `monteCarlo()` function. We advice the users to only use this model for small sequences **with only one or two hundreds components**. Further developpments will be made for markovian model in the next version of the package.

```
MyTransMat <-
matrix(c(0.3,0.1,0.1,0.1,0.4, 0.2,0.2,0.1,0.2,0.3, 0.3,0.4,0.1,0.1,0.1, 0.3,0.3,0.1,0.0,0.3,
        0.1,0.1,0.2,0.3,0.3), ncol = 5, byrow=TRUE)

MySeq.CM=transmatrix2sequence(matrix = MyTransMat,length=150, score =-2:2)
MySeq.CM
#> [1] 2 2 -2 1 0 -2 0 -1 0 -2 2 1 -2 2 -2 -2 2 1 2 -1 2 0 -2 2 2
#> [26] 1 2 -1 2 2 0 -1 2 1 -2 1 -1 1 -2 2 1 -2 1 0 2 2 2 -1 2 2
#> [51] 2 1 -1 2 -1 -2 2 2 2 0 0 -2 2 1 2 0 1 -2 -1 -2 -2 -2 -1 2
#> [76] 1 0 -1 2 2 2 1 0 1 2 2 2 2 2 2 1 2 2 1 -1 2 2 0 -1 -2
#> [101] 0 -2 2 -1 0 -1 2 2 0 -2 0 -1 -2 2 2 -1 -2 -2 0 -1 -1 -1 -2 -1 -1
#> [126] 1 -2 2 1 -1 2 1 -1 -2 2 0 -2 2 1 -1 -1 2 2 2 1 0 -2 2 2 1
AA.CM=automatic_analysis(sequences = list("x1" = MySeq.CM), model = "markov")
#>
#> |
#> |
#> |
#> |=====| 100%
AA.CM
#> $x1
#> $x1$p-value`
#> [1] 0.4693897
#>
#> $x1$method applied`
#> [1] "Exact Method"
#>
#> $x1$localScore
#> $x1$localScore$localScore
```

```

#> value begin end
#> 54 17 97
#>
#> $x1$localScore$suboptimalSegmentScores
#> value begin end
#> [1,] 4 1 2
#> [2,] 3 11 12
#> [3,] 54 17 97
#>
#> $x1$localScore$RecordTime
#> [1] 1 11 17

```

With MonteCarlo method the local score p -value is also not significant but not very accurate. Number of simulation can not be too large as it induces a increased time computation.

```

Ls.CM=AA.CM$x1$localScore[[1]][1]
monteCarlo(local_score = Ls.CM,
FUN = transmatrix2sequence, matrix = MyTransMat,
length=150, score = -2:2,
plot=FALSE, numSim = 10000)
#> p-value
#> 0.0763

```