



**HAL**  
open science

# localScore: an R package to highlight optimal and suboptimal segment in a sequence with associated p-values computation

David Robelin, Sébastien Déjean, Sabine Mercier

## ► To cite this version:

David Robelin, Sébastien Déjean, Sabine Mercier. localScore: an R package to highlight optimal and suboptimal segment in a sequence with associated p-values computation. 2024. hal-04723307

**HAL Id: hal-04723307**

**<https://hal.science/hal-04723307v1>**

Preprint submitted on 7 Oct 2024

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# 1 localScore: an R package to highlight optimal and suboptimal segment in a sequence with associated $p$ -values computation

David Robelin<sup>1</sup>, Sébastien Déjean<sup>2</sup>, and Sabine Mercier<sup>3</sup>

<https://doi.org/10.5802/fake.doi>

## Abstract

Highlighting atypical segments of a sequence is an important goal in very diverse domains. In the case where no prior information on the length of the segment to be highlighted is known, Karlin and Altschul defined in 1990 the local score for biological sequence analysis, and an asymptotic approximation of its distribution is proposed in 1992. There exist now many other theoretical results to establish the local score  $p$ -value in different contexts.

We developed an R package gathering these results for a sequence modeled by independent and identically distributed variables. It allows to compute the local score, the suboptimal scores, their position, and proposes to establish the local score  $p$ -value using the different theoretical methods available so far. An automatic analysis is also proposed to perform the most appropriate method according to the analyzed sequence. We present here the package and different examples of application. Comparisons with other tools used depending on the context of application are also given. The `localScore` package is available on the Comprehensive R Archive Network. It is distributed under the GPL-2 licence for the core program (and various licenses for embedded Eigen library)

**Keywords:** statistical significance, local score, sequence analysis, atypical segment detection, Lindley-CUSUM process

<sup>1</sup>GenPhySE, UMR1388, INRAE Toulouse, 24, chemin de Borde-Rouge, Auzeville Tolosane, France,

<sup>2</sup>Institut de Mathématiques de Toulouse, UMR 5219 Université de Toulouse et CNRS, Université Toulouse III Paul Sabatier, 118 route de Narbonne, Toulouse, France, <sup>3</sup>Institut de Mathématiques de Toulouse, UMR 5219 Université de Toulouse et CNRS, Université Toulouse II Jean Jaurès

## Correspondence

[machin@example.edu](mailto:machin@example.edu)

## Introduction

2  
3 Highlighting atypical periods or segments in sequences is an issue of interest in many fields,  
4 such as Bioinformatics and Genomics, Biosurveillance, Ecology and Environmental Sciences, Epi-  
5 demiology and Health Sciences, Finance, Reliability and Quality Control, Telecommunication Sci-  
6 ences, and many others... Karlin and Altschul, 1990 defined the local score statistic to analyze  
7 biological sequences: it corresponds to the maximum cumulative value of a given property over  
8 every possible segments in a sequence, considering segments of any position and any length (see  
9 Equation (1)). The calculation of the statistical significance of the local score is crucial in order  
10 to distinguish atypical values from ones that could have appeared by chance. Karlin and Dembo,  
11 1992 proposed asymptotic approximations of the distribution of the local score when the length  
12 of the sequence is growing to the infinity. A generalization of this approximation for the sequence  
13 comparison case is developed in BLAST Software<sup>1</sup>, but to our knowledge no development have  
14 been done for a single sequence analysis case. At the present time, results exist that consider in-  
15 dependent or dependent models on the sequence. Those results include: improvements of the  
16 approximations of Karlin *et al.* (see Cellier et al., 2003 for the independent model and Grusea  
17 and Mercier, 2020 for the Markov model); exact methods (see Mercier and Daudin, 2001 for  
18 the independent model and Hassenforder and Mercier, 2007 for the Markov model); a result on  
19 the distribution for the pair of the local score value and the length of the segment that realizes  
20 the local score value (see Chabriac et al., 2014; Lagnoux et al., 2017). We developed the package  
21 `localScore` (Simon et al., 2023) for the software R (R Core Team, 2024). In this first version of  
22 the package, we focus on the different ways to establish the statistical distribution of the local  
23 score in a sequence modeled by independent and identically distributed (I.I.D.) random variables,  
24 and address a first result in the Markovian case.

25 The remainder of the article contains a brief presentation of the main theoretical backgrounds  
26 implemented in the `localScore` package. Then the package is described in Section through a  
27 standard workflow to follow and a description of the main functions. Section presents examples  
28 of using `localScore` in four different domains: biological sequence analysis and a comparison  
29 with a sliding window statistics; a signal detection context and a comparison with control charts;  
30 epidemiology and a comparison with scan statistics; a genomic sequence analysis.

## Theoretical background

31  
32 Let us consider a sequence as a succession of components that belong to a finite set  $\mathcal{A}$ . It  
33 can be for example a DNA sequences with  $\mathcal{A} = \{A, C, G, T\}$ . Let us define a score scheme  
34 or a score scale as a function  $s$  that assigns a real number to any letter of  $\mathcal{A}$ . The score can, for  
35 example, quantify a physico-chemical property. See the web site of Protscale<sup>2</sup> for an illustration  
36 of different score scales in biological sequence analysis context. Let  $\mathbb{A} = (A_i)_{1 \leq i \leq n}$  be a sequence,  
37 and let us denote  $X_i := s(A_i)$ , for  $i \geq 1$ , the scoring sequence associated to the sequence  $\mathbb{A}$  based  
38 on the score function  $s$ . Examples of scoring functions are presented in Section .

<sup>1</sup><https://blast.ncbi.nlm.nih.gov/>

<sup>2</sup><https://web.expasy.org/protscale>

39 With  $X_0 := 0$ , the local score  $M_n$  of one sequence  $\mathbb{X}$  of length  $n$  is defined by

$$(1) \quad M_n := \max_{0 \leq i \leq j \leq n} \sum_{k=i}^j X_k .$$

40 In Mercier and Daudin, 2001 the authors proved that the local score can also be defined as:  
 41  $M_n := \max_{0 \leq i \leq n} U_i$  with  $U_0 := 0$  and  $U_{i+1} := \max(U_i + X_{i+1}, 0)$  the Lindley, or CUSUM process,  
 42 associated to the sequence  $(X_i)_{1 \leq i \leq n}$ . The Lindley process defines non negative excursions and  
 43 the height of the highest one is equal to the local score. The other excursions are called the sub  
 44 optimal segments.

45 The local score approach avoids the choice of a segment length when no prior information  
 46 on it is available. Let us present below the two main kinds of results, approximation and the exact  
 47 method, when the random variables  $(A_i)_{1 \leq i \leq n}$  are independent and identically distributed (I.I.D.)  
 48 and so are the  $(X_i)_{1 \leq i \leq n}$ .

#### 49 Karlin and Dembo approximation

50 The asymptotic approximation of Karlin and Altschul, 1990 and Karlin and Dembo, 1992 cor-  
 51 responds to an asymptotic result converging to a Gumbel distribution when the length of the  
 52 sequence  $n$  tends to infinity. This result stands on the two following hypotheses: The average  
 53 score must be non positive,  $\mathbb{E}[X] < 0$ , and a non negative score must be possible,  $\mathbb{P}(X > 0) > 0$ .  
 54 We have

$$(2) \quad \lim_{n \rightarrow +\infty} P \left( M_n \leq \frac{\ln n}{\lambda} + x \right) = e^{-K^* e^{-\lambda x}}$$

55 where  $\lambda$  and  $K^*$  depend on the score distribution. The parameter  $\lambda$  corresponds to the single  
 56 root of a polynomial of degree equal to the amplitude of the scores (maximum score minus  
 57 minimum score) and checking  $\mathbb{E}[\exp(\lambda X)] = 1$ . The existence of  $\lambda$  is ensured by the assumption  
 58  $\mathbb{E}[X] < 0$ . The set of other roots is also used for the calculation of  $K^*$  notably by means of a  
 59 square matrix called Vandermonde comprising at each line a geometric progression associated  
 60 with one of the roots of the polynomial. This gives the following approximation for an observed  
 61 local score  $a$

$$(3) \quad P(M_n \leq a) \approx e^{-K^* n e^{-\lambda a}} .$$

62 The computation of the approximation given in (3) is very accurate for sequence length larger  
 63 than thousands and very fast to obtain but must be avoided for sequence shorter than a hundred  
 64 components.

#### 65 Karlin by Monte Carlo

66 The Karlin and Dembo approximation in (2) calculates the value of two parameters  $\lambda$  and  $K^*$   
 67 in function of the values and distribution of the scores. Here we propose to estimate them by a  
 68 Monte Carlo approach. This method is useful in the case of a too long score sequence to perform  
 69 a direct and efficient Monte Carlo of the local score distribution as it does not need to simulate

70 full length sequences. Formula (3) can be linearized in  $\lambda$  and  $K^*$  using logarithms as long as  $n$  is  
71 large enough. That leads to the following formula:

$$(4) \quad \ln \{-\ln \{P(M_n \leq a)\}\} \approx \ln K^* - \lambda a + \ln n .$$

72 Given the previous formula, the Karlin by Monte Carlo procedure consists in:

- 73 (1) Choosing a sequence length  $n_{sim}$  for the simulation big enough to have a satisfying Karlin  
74 and Dembo approximation and small enough to be computed with reasonable resources.
- 75 (2) Simulating sequences of size  $n_{sim}$ .
- 76 (3) Calculating the local score of each sequence in order to derive empirical distribution  
77 function of the local score for sequence of size  $n_{sim}$ .
- 78 (4) Deriving estimation of  $\lambda$  and  $K^*$  by a linear regression on the empirical distribution func-  
79 tion using Formula (4), i.e.,  $\hat{\lambda} = -\hat{b}$  and  $\hat{K}^* = \exp(\hat{a})/n_{sim}$  where  $\hat{a}$  and  $\hat{b}$  are respectively  
80 the slope and the intercept of the regression.
- 81 (5) Using Karlin and Dembo approximation to calculate the  $p$ -value of the local score ob-  
82 served on the full sequence of size  $n$ .

### 83 Daudin

84 For  $a$  an observed local score value, the exact method in the I.I.D. case is based on an appro-  
85 priate stopped process constructed to be a Markov chain and taking its values in  $\{0, \dots, a\}$ . Let  
86 us denote  $P = (P_{ij})_{0 \leq i, j \leq a}$  its corresponding transition matrix. Mercier and Daudin, 2001 proved  
87 that

$$(5) \quad (\forall a \geq 0) \quad \mathbb{P}(M_n \geq a) = (P^n)_{(0,a)} .$$

88 There is no restriction on the sign of the average score for the use of the exact method. This  
89 method is accurate and very fast for  $n$  up to several thousands but must be avoided for very long  
90 sequence of order a million because it could be too time and space consuming. As there exists  
91 a limit to the exponentiation of  $P$  for  $n$  tending to infinity, it is not necessary to use the correct  
92 value of  $n$  for sequence of length larger than a hundred of thousands, to obtain an accurate value,  
93 but a smaller value could be used.

### 94 Improved approximation of Karlin et al. in the I.I.D. case

95 An improved approximation of the one proposed in Karlin and Dembo, 1992 is proposed in  
96 Cellier et al., 2003. As with the Karlin *et al.* method, it is necessary to calculate the roots of the  
97 same polynomial which are then used in several steps in order to calculate the additive correcting  
98 terms to improve the approximation of Karlin *et al.* We have for large  $a$  values

$$(6) \quad P(M_n \leq a) \approx (1 - \sum_{i=1}^{\mu} K_i R_i^a)^{\frac{n}{\mu} + 1}$$

99 with  $(R_i)_i$  the roots of module strictly less than 1 of a polynomial directly defined with the  
100 score distribution. The degree of this polynomial is equal to the range of the possible scores.  
101 Based on the two hypothesis used in the work of Karlin and Dembo, there exists a unique positive  
102 real root with module less than 1,  $e^{-\lambda}$ , with  $\lambda$  defined in Equation (2). The parameters  $K_i$  and  
103  $\mu$  are also derived from the score distribution and the computation based on the Vandermonde

104 matrix and some equation system resolutions. The improved approximation in (6) is accurate and  
 105 fast for values of  $n$  from several hundreds, but must be avoided for sequence length lower than  
 106 one hundred.

107 All the above theoretical results must be considered complementary for practical application  
 108 depending on the score scheme, with its range, the sign of the average score, and the length of  
 109 the sequence to be analyzed.

## 110 Software features and contents

### 111 Workflow

112 A tentative workflow using `localScore` could be:

- 113 (1) Transform the component of a given sequence set into score sequences through a given  
 114 score function.
- 115 (2) Learn the distribution of the scores on the score sequences.
- 116 (3) Compute the local score of each sequence.
- 117 (4) Compute the corresponding  $p$ -values using the automatic method for the computed local  
 118 score value, the corresponding sequence length and the global score distribution.

### 119 Main functions

120 Following the workflow presented above, here are the main functions that can be used in  
 121 each step.

122 *To get a score sequence:* The transformation of a component sequence, as a DNA one, into a  
 123 score sequence can be done with the `CharSequence2ScoreSequence` function. Integer or  
 124 real scores can be considered.

125 *To learn a distribution:* Learning distribution of the components of the sequences or the given  
 126 scores can be performed using several functions. For instance, the empirical distribution from  
 127 one numerical sequence or a list of sequences is built by `scoreSequences2probabilityVector`.

128 *To compute the local score:* The function `localScoreC` (respectively `localScoreC_double`)  
 129 calculates the local score for a sequence of integer (*resp.* real) scores. It provides the local score  
 130 and all suboptimal segments with associated scores. Functions `suboptimalSegment` or `Lindley`  
 131 can be used to obtain the others localizations of the different realizations of the local score.

132 *To compute the corresponding  $p$ -values:* Then, the following functions propose different meth-  
 133 ods to compute  $p$ -values associated to the local score of a sequence:

- 134 • `karlin`: The Karlin *et al*'s approximation (see (3)). This method needs a non positive  
 135 average score,  $\mathbb{E}[X] < 0$ , and integer scores, and is more adapted for long sequences  
 136 with length larger than a few thousand components, depending on the expectation of  
 137 the score distribution.
- 138 • `mcc`: An improved approximation of the previous one presented in Cellier et al., 2003.  
 139 This method also needs a non positive average score,  $\mathbb{E}[X] < 0$ , and integer scores, and  
 140 is more adapted for sequences with length from a few hundreds components, depending  
 141 on the expectation of the score distribution.
- 142 • `daudin`: An exact method for integer scores is also incorporated and can be used what-  
 143 ever the sign of the expected score (see (5)). This method is computationally adapted  
 144 for not too long sequences, but several thousands of components can be easily handled.

**Table 1** – Adequate methods to compute the local score  $p$ -value depending on the average score value  $\mathbb{E}[X]$  and the sequence length  $n$  order ; with E : `daudin()` ; MCC : `mcc()` ; K : `karlin()` ; MC : `monteCarlo()` ; MC-K : `karlinMonteCarlo()`.

$n$	$< 100$	$10^2 \leq \cdot < 10^3$	$10^3 \leq \cdot < 10^4$	$\geq 10^4$
$\mathbb{E}[X] < 0$	E ; MC	E ; MCC ; MC	E ; MCC ; MC	MCC ; K ; MC ; MC-K
$\mathbb{E}[X] \geq 0$	E ; MC	E ; MC	E	

145 The implementation is based on the exponentiation of a square matrix of size  $a$ , with  $a$  a  
146 given local score value.

- 147 • `monteCarlo`: A classical Monte Carlo method
- 148 • `karlinMonteCarlo` and `karlinMonteCarlo_double`: A mix between the Karlin *et al.*'s and the Monte Carlo method. It allows an approximated distribution with a lower  
149 time computation than the empirical Monte Carlo method, for very long sequences. This  
150 mixed method also needs  $\mathbb{E}[X] < 0$ .

152 We also developed the function `automatic_analysis` for users with less experience. This  
153 function, as its name indicates, automatically picks the adequate  $p$ -value method for the user's in-  
154 put according to the configuration described in Table 1. The function calculates the  $p$ -value based  
155 on the length of each of the sequences given as input. It can either use an empirical score distribu-  
156 tion based on the input or a distribution provided by the user. By setting the `method_limit`, the  
157 user can also decides up to what sequence length the computation-intensive methods (`daudin`,  
158 `exact_mc`) should be used to calculate the  $p$ -value.

### 159 Inputs / outputs

160 *Inputs.* When starting the workflow, the first input is a sequence. It can be imported in R from  
161 an ASCII file using the standard reading functions such as `read.table` and related functions.  
162 For users interested in analyzing biological sequences composed of nucleotides or amino acids,  
163 the package can also handle FASTA files as inputs. In FASTA files, every sequence is preceded by  
164 a title (marked by a ">") and a line break. One sequence takes one line, followed by a line break  
165 and a line only containing a tab.

166 Furthermore, if no sequences are passed to the `automatic_analysis` function, it let the  
167 user pick a FASTA file. In this case, and if the user hasn't provided any score system (as it can  
168 be done by passing a named list with the appropriate scores for each character), the second file  
169 dialog pops up. The latter allows to choose a file containing the score, and if the user provides an  
170 extra column for the probabilities, they are used, too - see Section File Formats in the vignette  
171 for details.

172 Score files can also be imported in a standard way from an ASCII file. Such a file must contain  
173 a header and each row contains a letter and its score. Optionally, a probability for each score  
174 can also be provided.

175 *Numerical outputs.* The main numerical output is given by the `localScoreC` function. It con-  
176 tains a list with the following attributes:

- 177 • The local score value and the begin and end index of the segment realizing this optimal  
178 score.
- 179 • All the local maxima of the Lindley process (non negative excursion) and their begin and  
180 end index.

- 181       • The record times of the Lindley process but only the ones corresponding to the begin  
182       index of non negative excursions.

183       Every method calculating  $p$ -values only provide the value obtained.

184 *Graphical outputs.* Graphical outputs can be optionally displayed by the `monteCarlo` and the  
185 `karlinMonteCarlo` functions. They represent the distribution of all local scores simulated and  
186 the cumulative distribution.

### 187 **Example data**

188       Some data we propose to analyze in Section are already embedded in the package for illustra-  
189 tion purpose. `Seq1093` is a real biological sequence with 1093 characters referring to Q60519  
190 queries in UniProt Data base<sup>1</sup>. `SeqListSCOPE` contains 285 protein sequences with length  
191 from 31 to 404. They are referred as `CF_scop2dom_20140205aa` in the Structural Classifica-  
192 tion Of Proteins database (SCOP)<sup>2</sup>. `SJSyndrome.data` corresponds to a dataset of 824 lines,  
193 each describing a Stevens-Johnson syndrome appearance described by 15 covariates including  
194 Case ID, Initial FDA Received Date, days since last fda. The third column corresponds to the  
195 number of days between two adverse events. `Aeso` consists of individual dates of birth over 35  
196 cases of the birth defects oesophageal and tracheo-oesophagean fistula observed in a hospital  
197 in Birmingham.

## 198 **Illustrations**

199       We illustrate the use of the `localScore` package on four examples in different fields. First,  
200 one of the biological sequences embedded in the package is used as a toy example to show  
201 a basic use of the package. In the same vein, we illustrate how to deal simultaneously with  
202 a set of sequences. Then, we analyze two medical data sets to show how local score can be  
203 used to detect eventual shift in sequential observations. Subsection deals with the study of a  
204 chromosome to associate genomic regions with phenotype differentiation. We also present, for  
205 each case, results of other methods.

### 206 **Biological sequences**

207       We first describe how to analyze one single sequence then we show how to deal with a set  
208 of several sequences at once.

209 *One single sequence.* Several sequences are already embedded in the package. Let us use the  
210 `Seq1093` object, corresponding to the protein Q60519 SEM5B\_MOUSE<sup>3</sup>. With 1093 charac-  
211 ters, we consider it as a sequence for which quite all the possible proposed methods can be used  
212 to establish the statistical significance (see Table 1).

```
213 R> library(localScore)
214 R> data(Seq1093)
215 R> MySeq <- Seq1093
216 R> nchar(MySeq)
```

<sup>1</sup><https://www.uniprot.org>

<sup>2</sup><https://scop.mrc-lmb.cam.ac.uk/>

<sup>3</sup><https://www.uniprot.org/uniprot/Q60519>



217 The function `CharSequence2ScoreSequence` converts the character sequence into a  
 218 score sequence using the `HydroScore` object providing the correspondence between a letter  
 219 and its score according to Kyte & Doolittle hydrophobic score scale (Kyte and Doolittle, 1982).

```
220 R> data(HydroScore)
221 R> SeqScore <- CharSequence2ScoreSequence(MySeq, HydroScore)
```

222 Then the local score computation can be performed.

```
223 R> ResLocalScoreMySeq <- localScoreC(SeqScore)
224 R> ResLocalScoreMySeq
```

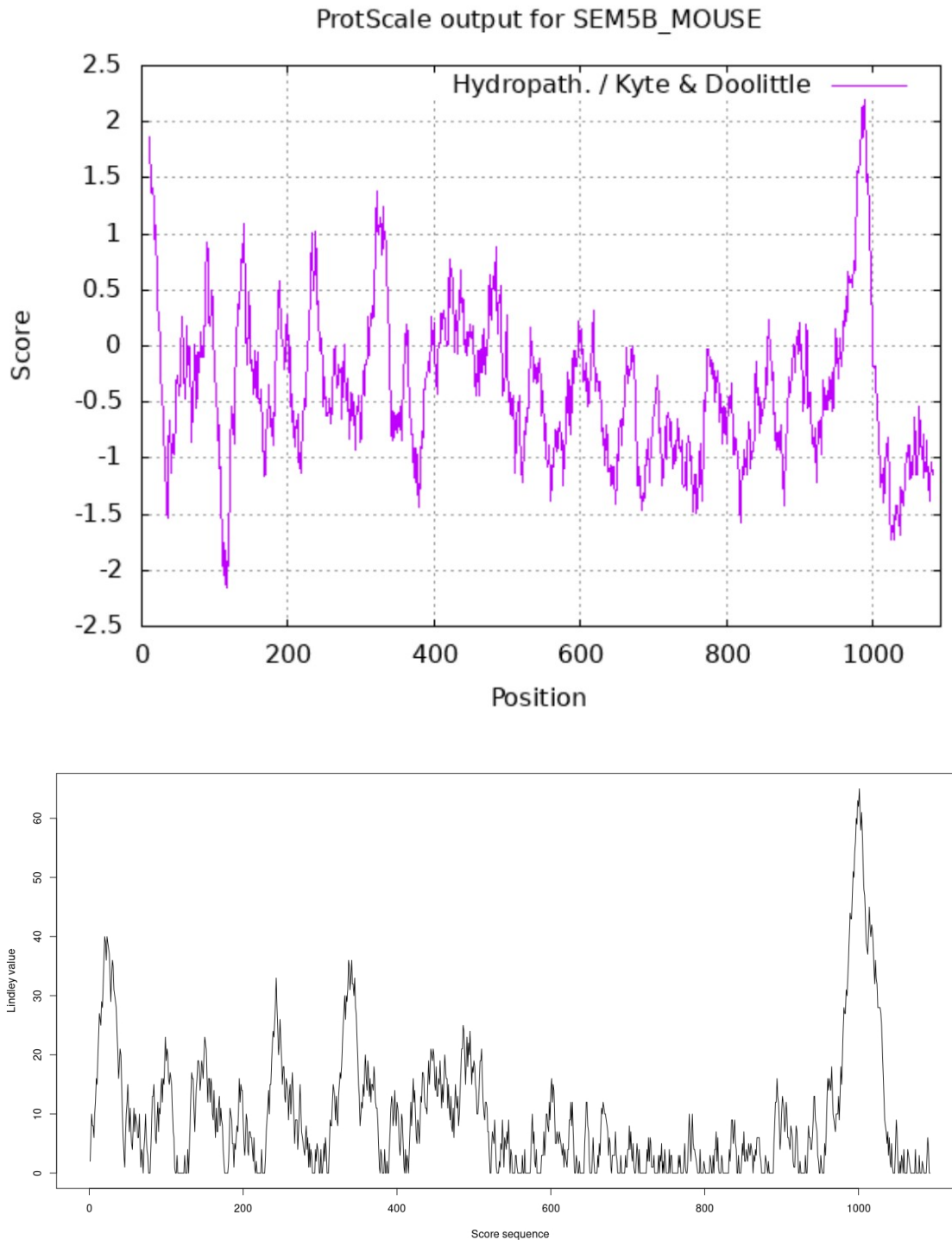
```
225 $localScore
226 value begin end
227 65 956 1001
228
229 $suboptimalSegmentScores
230 value begin end
231 [1,] 40 1 20
232 [2,] 10 71 73
233 [3,] 23 80 99
234 [4,] 3 114 114
235 [5,] 3 124 124
236 [6,] 4 128 128
237 [7,] 23 130 150
238 [8,] 16 181 195
239 [9,] 2 217 217
240 [10,] 4 224 224
241 [...]
242 [70,] 2 1054 1054
243 [71,] 3 1056 1056
244 [72,] 4 1059 1059
245 [73,] 4 1064 1064
246 [74,] 4 1074 1074
247 [75,] 3 1079 1079
248 [76,] 2 1083 1083
249 [77,] 6 1089 1090
250
251 $RecordTime
252 [1] 1 71 80 114 124 128 130 181 217 224 229
253 [12] 294 301 304 312 379 384 387 390 411 413 416
254 [23] 523 533 548 554 563 566 574 588 633 637 644
255 [34] 654 661 680 691 697 711 715 725 740 742 746
256 [45] 748 754 756 766 773 778 794 799 807 811 813
257 [56] 822 832 844 848 864 879 883 890 924 931 934
258 [67] 951 956 1048 1054 1056 1059 1064 1074 1079 1083 1089
```

259 We retrieve only the local score value for further use when calculating the  $p$ -value.

```

260 R> LocalScoreMySeq <- ResLocalScoreMySeq$localScore[1]
261     The function scoreSequences2probabilityVector builds an empirical distribution from
262 the sequence.
263 R> ProbDistribution <- scoreSequences2probabilityVector(SeqScore)
264 R> round(ProbDistribution, 3)
265     -5     -4     -3     -2     -1     0     1     2     3     4     5
266 0.074 0.203 0.020 0.075 0.212 0.078 0.000 0.071 0.094 0.144 0.028
267     The exact method (see (5)) can then be used to compute the  $p$ -value.
268 R> ResDaudin <- daudin(localScore = LocalScoreMySeq,
269 +   sequence_length = length(SeqScore),
270 +   score_probabilities = ProbDistribution,
271 +   sequence_min = min(SeqScore),
272 +   sequence_max = max(SeqScore))
273 R> ResDaudin
274 [1] 0.072
275     The approximate method of Karlin et al. (see (3)) can be performed equivalently with the
276 karlin function.
277 R> ResKarlin <- karlin(localScore = LocalScoreMySeq,
278 +   sequence_length = length(SeqScore),
279 +   score_probabilities = ProbDistribution,
280 +   sequence_min = min(SeqScore),
281 +   sequence_max = max(SeqScore))
282 R> ResKarlin
283 [1] 0.076
284     The two  $p$ -values are rather close (0.072 for the exact method, 0.076 for the approximate
285 one).
286     In comparison, here are the results obtained with ProtScale Expsy web tool on the same
287 sequence. ProtScale computes and represents the profile on a selected protein produced by
288 any amino acid scale and accumulating the score values over a sliding window of a chosen size.
289 Note that the possible window size are restricted to odd values from 3 to 21. We used the
290 hydropathicity scale proposed by Kyte and Doolittle, 1982. We chose a size equal to 21 which is
291 the closest value to the length of the optimal segment given by the local score approach without
292 any prior information on it. The results are presented in Figure 1. We can observe one main peak
293 and the numerical output (not shown) gives us a window value equal to 2.195 and a center index
294 equal to 989 (begin index 979; end index 999). This segment corresponds to the one highlighted
295 by the local score but with a length equal to 45 with begin index 956 and end index 1001. The
296 local score  $p$ -value allows us to say that this region is not statistically significant. For a window
297 size equal 9 corresponding to the one given by default, we can observe several picks with a
298 similar value before the one we discuss previously. We have also represented the corresponding
299 Lindley process using the {lindley} function.
300 R> LindleySeqScore <- lindley(SeqScore)
301 R> plot(LindleySeqScore, type="l")

```



**Figure 1** – Top: Graphical output of the results provided by the ExPASy ProtScale web tool for the corresponding sequence Q60519, the Kyte and Doolittle scale and a window size equal to 21. Bottom: Lindley process calculated with the `localScore` package.

302 *A set of sequences.* The data consists in a list of 285 character strings with their entry codes as  
 303 names extracted from the Structural Classification Of Proteins database (SCOP)<sup>2</sup>. More precisely  
 304 this data contain the 285 protein sequences of the data called “CF\_scop2dom\_20140205aa”  
 305 with sequence length from 31 to 404.

<sup>2</sup><https://scop.mrc-lmb.cam.ac.uk/>

306 This sequence is a part of the package and can be loaded and briefly explored with:

```
307 R> data(SeqListSCOPE)
308 R> SeqListSCOPE[1]
309 P50456
310 "ARDVIQVVIDHNVGAGVITDGHLLHAGSSSLVEIGHTQVDPYGKRCYCGNHGCLLETIAS
311 VDSILELAQLRLNQSMSSMLHGQPLTVDSLQQAALRGDLLAKDIITGVGAHVGRILAIMV
312 NLFNPQKILIGSPLSKAADILFPVISDSIRQQALPAYSQHISVEST"
```

```
313 R> nchar(SeqListSCOPE[1])
```

```
314 P50456
```

```
315 165
```

```
316 R> summary(sapply(SeqListSCOPE, nchar))
```

```
317   Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
```

```
318   31.0    78.0   102.0   121.7   141.0   404.0
```

319 The sequence lengths varie from 31 to 404.

320 The function `CharSequence2ScoreSequence` transforms the protein sequence into a  
321 score sequence using the `HydroScore` object. The score corresponding to each amino acid  
322 can be displayed as:

```
323 R> data(HydroScore)
```

```
324 R> unlist(HydroScore)
```

325 and the conversion is done as follows:

```
326 R> MySeqScoreList <- lapply(SeqListSCOPE,
```

```
327 +   FUN = CharSequence2ScoreSequence, HydroScore)
```

328 Then we use `automatic_analysis` function to perform the most appropriate method to  
329 compute the *p*-value of the local score of each sequence.

```
330 R> ResAutoAnalysis <- automatic_analysis(sequences = MySeqScoreList,
```

```
331 +   model='iid')
```

332 The results can then be investigated.

```
333 R> ResAutoAnalysis[[1]]
```

```
334 $`p-value`
```

```
335 [1] 0.06389172
```

```
336
```

```
337 $`method applied`
```

```
338 [1] "Exact Method Daudin et al"
```

```
339
```

```
340 $localScore
```

```
341 $localScore$localScore
```

```
342 value begin end
```

```
343 67 4 144
```

```
344
```

```
345 $localScore$suboptimalSegmentScores
```

```
346 value begin end
```

```
347 [1,] 2 1 1
```

```

348 [2,]      67      4 144
349
350 $localScore$RecordTime
351 [1] 1 4

```

352 The first sequence of the list has a local score value equal to 62 and the segment that realizes  
353 this maximum begins at index 4 and finishes at index 144. Its  $p$ -value equals 6.39%.

354 We can easily extract the first 10  $p$ -values, the 5 smallest  $p$ -values, the significant sequences  
355 and their local score values.

```

356 R> sapply(ResAutoAnalysis, function(x){x$`p-value`})[1:10]
357 P50456 P14859 P10037 Q13619 P22262 P20823 P07014 Q9X399 Q0SB06 Q9I641
358 0.064 0.973 0.875 0.896 0.451 0.967 0.749 0.681 0.994 0.512
359 R> sort(sapply(ResAutoAnalysis, function(x){x$`p-value`})) [1:5]
360      Q5SMG8      P0A334      Q2W6R1      O27564      P12282
361 9.485100e-07 3.442818e-04 4.406208e-04 4.548065e-04 6.167591e-02
362 R> which(sapply(ResAutoAnalysis, function(x){x$`p-value`}) < 0.05)
363 Q2W6R1 O27564 P0A334 Q5SMG8
364      14      90      150      192

```

365 The local score of every sequence can be displayed as a boxplot or an histogram (Fig. 2):

```

366 R> SeqLocalScore <- sapply(ResAutoAnalysis,
367 +   function(x){x$localScore$localScore[1]})
368 R> boxplot(SeqLocalScore, horizontal = TRUE)
369 R> hist(SeqLocalScore)

```

370 The methods used to compute the  $p$ -values can be retrieved with

```

371 R> table(sapply(ResAutoAnalysis, function(x){x$`method`}))
372 Exact Method Daudin et al
373                               206

```

374 The maximum sequence length equals 404 so it is here normal that the exact method is used  
375 for all the 606 sequences of the data base. The score distribution that has been used to compute  
376 the  $p$ -value for every local scores is the empirical one estimated on the whole data set. It can be  
377 exhibited using

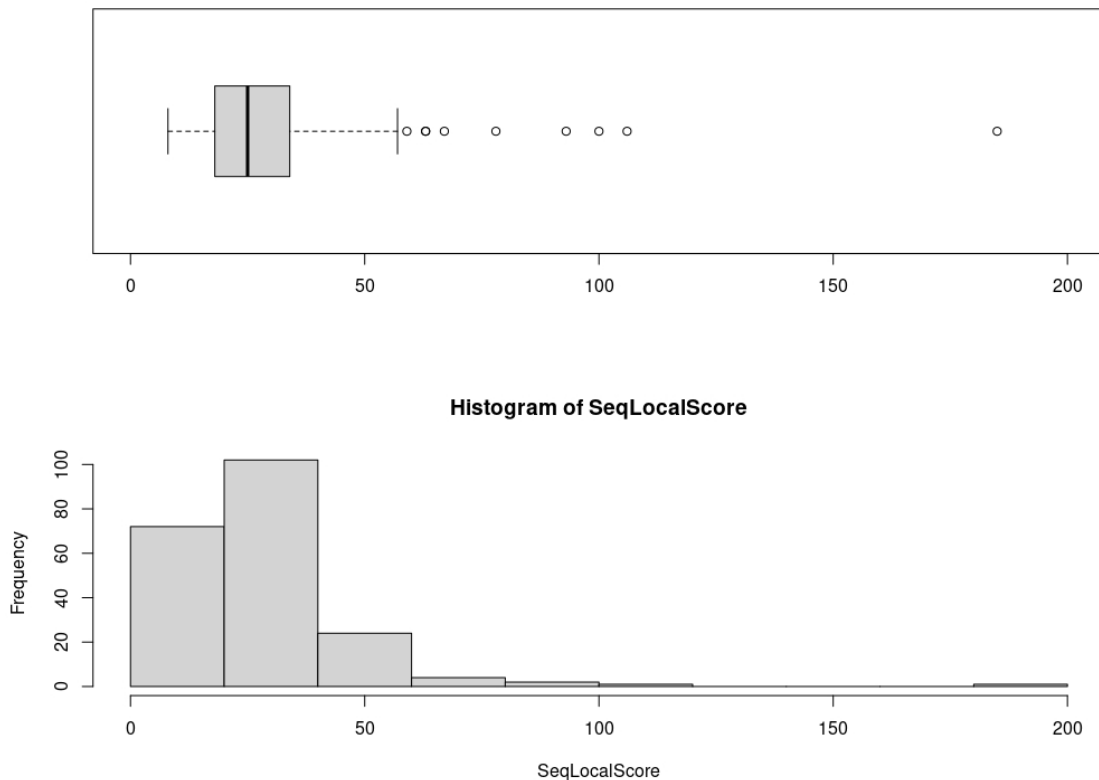
```

378 R> scoreSequences2probabilityVector(MySeqScoreList)
379      -5      -4      -3      -2      -1      0      1      2      3      4      5
380 0.055 0.264 0.022 0.041 0.148 0.072 0.000 0.105 0.052 0.175 0.067
381

```

## 382 Medical data

In other domains, as Telecommunication Sciences or Quality Control to name only those, where the goal is to highlight a change or a break point in the signal sequence, the data are analyzed as soon as there are collected. In such application domains, score scales are not previously proposed or constructed as it is done in biological sequence analysis. When testing at each time  $i$ , the null hypothesis  $H_0$ : "The observations  $(A_k)_{1 \leq k \leq i}$  follow the distribution  $f_\theta$  with parameter



**Figure 2** - Distribution of the local score of every sequence in the object `MySeqList`.

$\theta = \theta_0$ " vs  $H_1$ : "The observations  $(A_k)_{1 \leq k \leq i}$  follow  $f_{\theta_1}$  with  $\theta_1 \neq \theta_0$ ", it is usual to define the score of a given observation  $A_i$  at time  $i$  by the following Log Likelihood Ratio:

$$x_i = s(A_i) = \ln \left( \frac{f_{\theta_1}(A_i)}{f_{\theta_0}(A_i)} \right) .$$

383 Such a score function is used in this subsection and in the Subsection .

384 The local score can also be used to detect eventual shift in sequential observations. We  
 385 propose here to analyze data on the apparition of the Stevens-Johnson syndrome, a serious  
 386 dermatological disease due to a drug allergy.

```
387 R> data(SJSyndrome.data)
```

```
388 R> dim(SJSyndrome.data)
```

```
389 [1] 824 15
```

```
390 R> SJSyndrome.data[1:2,1:5]
```

```
391 Case.ID Initial.FDA.Received.Date days.since.last.fda Event.Date
```

```
392 1 4227848 10/01/1969 NA 02/16/1969
```

```
393 2 4227553 10/01/1969 0 07/10/1969
```

```
394 Latest.FDA.Received.Date
```

```
395 1 01-OCT-1969
```

```
396 2 01-OCT-1969
```

397 The third column `days.since.last.fda` corresponds to the number of days since the  
 398 last event (the Time Between Event sequence). The data present 824 adverse event apparitions  
 399 that lead to 823 Time Between two adverse Events (TBE) values in days.

```
400 R> DatesTBE <- SJSyndrome.data[-1,3] # the TBE sequence
401 R> n <- length(DatesTBE)
```

402 The TBE sequence can be modeled by a geometrical distribution. An estimation of its param-  
 403 eter is given by

```
404 R> p0Hat <- 1/(mean(DatesTBE[1:n]) - 1)
405 R> p0Hat
406 [1] 0.045349
```

407 with an estimated value equal to 0.045349 corresponding to the probability of observing  
 408 an adverse event at a given day among the whole studied population. Let us denote  $(T_i)_{1 \leq i \leq n}$   
 409 the TBE observations. At each time  $i$ , we want to test the following hypotheses :  $H_0$  "The obser-  
 410 vations  $(T_k)_{1 \leq k \leq i}$  follow a geometrical distribution with parameter  $p_0$ " vs  $H_1$  "The observations  
 411  $(T_k)_{1 \leq k \leq i}$  follow a geometrical distribution with parameter  $p_1 = 1.5 \cdot p_0$ ", with  $p_0$  and  $p_1$  in  $]0, 1[$ .

Let us define:

$$LLR(T) = \ln \frac{f_1(T)}{f_0(T)}$$

412 with  $f_j$  the probability density function of a geometrical variable of parameter  $p_j$  for  $j = 0, 1$ . At  
 413 each time  $i$  the local score of the sequence  $(LRR(T_k))_{1 \leq k \leq i}$  and its corresponding  $p$ -value are  
 414 computed using the package. More precisely, we compute  $LLR = \lfloor E \cdot \ln \frac{f_1(T)}{f_0(T)} \rfloor$  with  $E$  a tuning  
 415 parameter which allows a larger range of possible non negative scores. The use of this tuning  
 416 parameter does not change the segment that realizes the local score and neither its  $p$ -value (see  
 417 Fariello et al., 2017 Supplementary materials, for more details), but allows to highlight suboptimal  
 418 segments that could be interesting. We have here at least 3 non negative scores for  $E = 8$ .

```
419 R> p0 <- round(p0Hat, 4)
420 R> p1 <- 1.5*p0
421 R> E <- 8
```

422 Let us compute the score sequence and the local score for each sequence up to index  $i$  for a  
 423 sequential analysis.

```
424 R> ScoreSeq <- floor((log(dgeom(DatesTBE, p1) /
425 +   dgeom(DatesTBE, p0))) * E)
426 R> head(ScoreSeq)
427 [1] 3 -3 3 -15 -3 3
428 R> VectLS <- vector(length = n)
429 R> for (i in 1:n) {
430 +   VectLS[i] <- localScoreC(ScoreSeq[1:i])$localScore[1]}
431 R> head(VectLS)
432 R> tail(VectLS)
433 [1] 3 3 3 3 3 3
434 [1] 189 189 189 189 189 189
```

435 An alarm can be defined when the  $p$ -value of an observed local score value is less than a given  
 436 nominal level, usually 5% or 1%. In order to establish the  $p$ -value, the distribution of the scores

437 under the  $H_0$  hypothesis is needed. It is possible to established theoretically this distribution,  
 438 but in order to have a lighter presentation here, we empirically estimate the score distribution  
 439 on the data.

```
440 R> PkCal <- table(sort(ScoreSeq)) / length(ScoreSeq)
441 R> head(PkCal)
442 R> tail(PkCal)
443 [1] -109 -75 -61 -56 -54 -46
444 0.0012151 0.0012151 0.0012151 0.0036452 0.0012151 0.0012151
445 [1] -2 -1 0 1 2 3
446 0.036452 0.093560 0.117861 0.134872 0.227217 0.160389
```

447 We can notice that not all the possible values, between the minimum and the maximum score,  
 448 are present. The vector of the score distribution must be fulfilled.

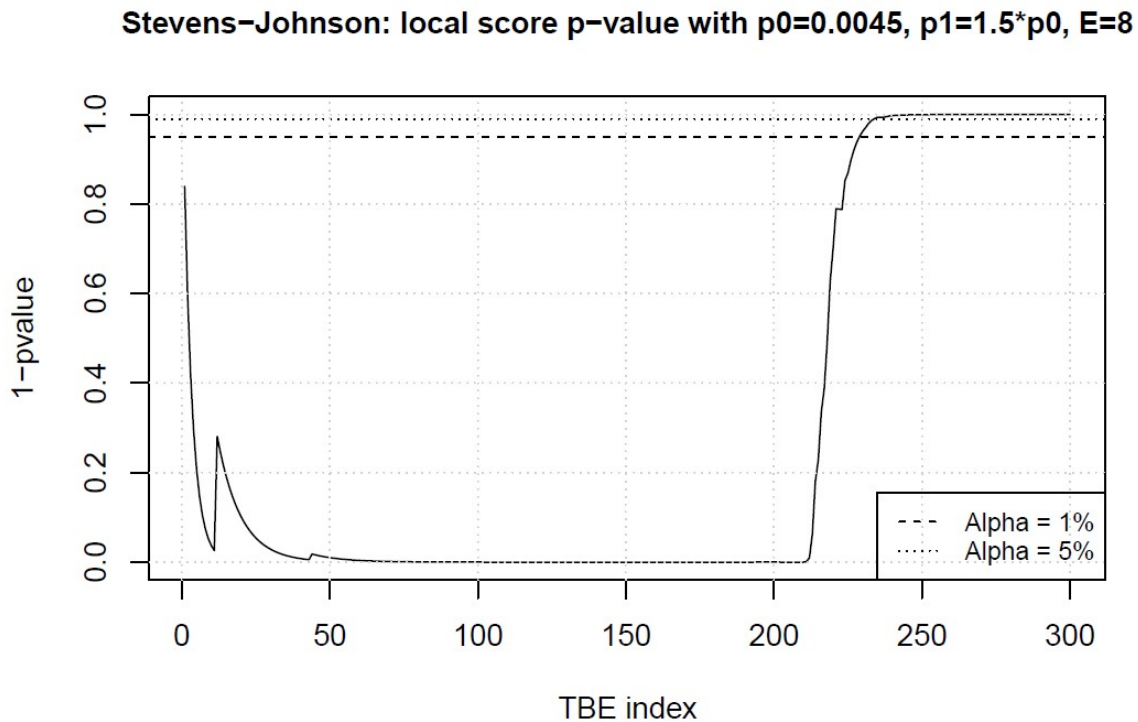
```
449 R> minXk <- min(ScoreSeq)
450 R> maxXk <- max(ScoreSeq)
451 R> score <- minXk:maxXk
452 R> PkEmp <- rep(0, length(score))
453 R> names(PkEmp) <- minXk:maxXk
454 R> for (i in 1:length(PkCal)) {
455 +   PkEmp[which(names(PkEmp) == names(PkCal)[i])] <- PkCal[i]}
456 R> head(PkEmp)
457 R> ProbaTh <- PkEmp
458 [1] -109 -108 -107 -106 -105 -104
459 0.0012151 0.0000000 0.0000000 0.0000000 0.0000000 0.0000000
460 R> mean(ScoreSeq)
461 [1] -1.633
```

462 The average score under  $H_0$  is non positive so any method, exact as approximated ones, can  
 463 be used to compute the statistical significance of the local score. We use the exact method with  
 464 `daudin` function as the sequence lengths allow it.

```
465 R> VectPval <- vector("numeric", length = n)
466 R> for (i in 1:n) {
467 +   LS <- localScoreC(ScoreSeq[1:i])$localScore[1]
468 +   VectPval[i] <- daudin(localScore = LS,
469 +     sequence_length = i,
470 +     score_probabilities = ProbaTh,
471 +     sequence_min = minXk,
472 +     sequence_max = maxXk) }
473 R> head(VectPval)
474 [1] 0.16039 0.40797 0.58157 0.70423 0.79093 0.85222
475 R> min(which(VectPval < 0.05))
476 [1] 229
```

477 Figure 3 illustrates the example on the first 300 observations where a first alarm, using a  
 478 nominal level  $\alpha = 5\%$ , appears at index 229.





**Figure 3** – Stevens Johnson syndrome: a unique alarm at index 229.

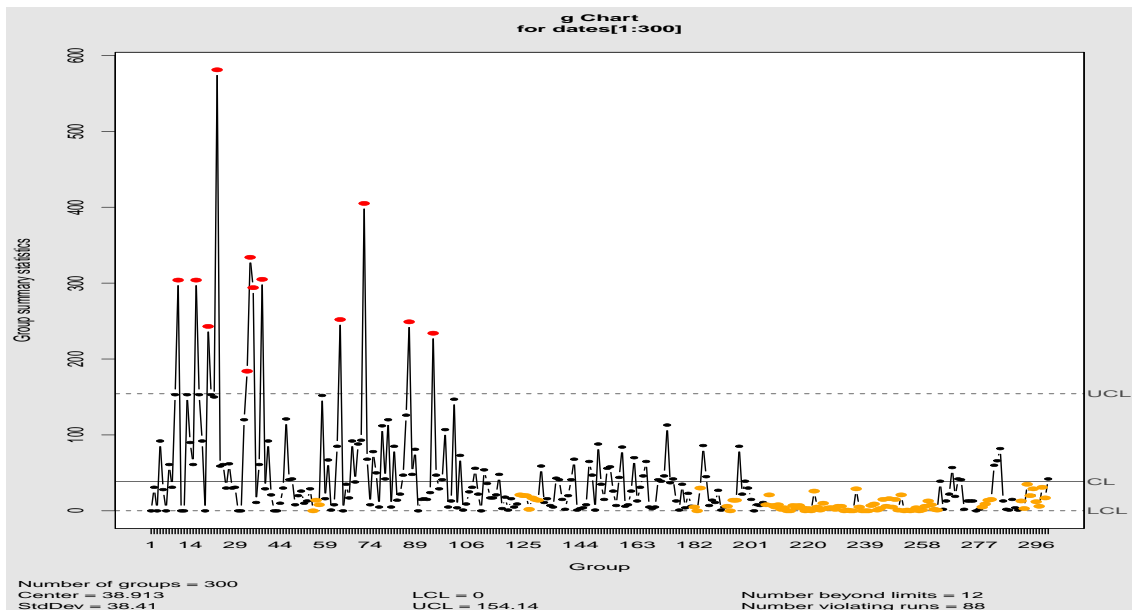
479 Different values for parameter  $p_1$  has been tested, and each case leads to a similar result.  
 480 Figure 3 representing the  $p$ -values at each index, can be seen as a control chart usually used  
 481 to analyze on-line sequences in industrial data (see for example the first and the most famous  
 482 control chart defined in 1930 and called the Shewhart chart, see W.A. Shewhart, 1931 mainly  
 483 used for Gaussian distribution). We can see one unique clear alarm at index 229. In Mercier,  
 484 2020 for a Gaussian model, it is shown that using the local score avoids false alarm better than  
 485 the usual control charts and allows to detect an existing change in the parameter in a competitive  
 486 mean time.

487 Let us see the Shewhart  $g$  chart, adapted for geometrical distribution, and proposed in the  
 488 package `qcc` in Figure 4.

```
489 R> library(qcc)
490 R> qcc(DatesTBE[1:300], type = "g")
491 R> Gchart <- qcc(DatesTBE[1:300], type = "g")
492 R> violating.runs(Gchart, run.length = qcc.options("run.length"))
```

493 Here the lower control limit (LCL) is equal to 0 and have no direct use. The twelve points up-  
 494 per than the upper control limit (UCL) in red are not “bad” alarm because they are corresponding  
 495 to a longer run than expected between two adverse events and are then considered as an im-  
 496 proved situation. We can observe several violating run in orange, corresponding to a particularly  
 497 numerous successive points under the central control limit, that are considering as alarms. One  
 498 is particularly long beginning at index 206 and including the index 229 of the alarm of the local  
 499 score chart.

500 Regarding both local score and  $g$  charts we suggest that the violating runs before index 206  
 501 in the  $g$  chart could be considered as false alarms.



**Figure 4** – Stevens Johnson syndrome - Shewhart  $g$  chart: A lot of alarms are pointed out. We can observe a violating run, corresponding to a particularly numerous successive points under the central control limit, beginning at index 206 and including the index 229 of the alarm of the local score chart.

## 502 Congenital oesophageal atresia data

503 The data consists of individual dates of birth over  $n = 35$  cases of the birth defects oe-  
 504 sophageal and tracheo-oesophagean fistula observed in a hospital in Birmingham, U.K., over  
 505 2191 days from 1950 through 1955, with Day one set as 1 January 1950 (see Knox, 1959). Glaz  
 506 et al., 2009 present in Chapter 17 different works on these data based on the use of scan statist-  
 507 tic. We first present in this section the results of the scan statistic analyses proposed in Glaz  
 508 et al., 2009 and secondly two different approaches based on the local score. The discrete scan  
 509 statistic  $S_{n,k}$ , with  $k \leq n$  two positive integers, of a sequence  $(S_i)_{1 \leq i \leq n}$  of  $n$  binary trials (1: suc-  
 510 cess, 0: failure) has been defined as the maximum number of successes within any  $k$  consecutive  
 511 trials. Let us consider a discrete sequence  $(S_i)_{i=1 \dots n}$ . We have  $S_{n,k} = \max_{1 \leq i \leq n-k+1} \sum_{j=i}^{i+k-1} S_j$ .  
 512 The data are given in the following line code. We also derive the Time Between two Events (or  
 513 success). The TBE sequence is modeled by a geometrical law of parameter  $p$ .

```
514 R> data(Aeso)
515 R> CasesIndex <- which(Aeso[,2] == 1)
516 R> tbe <- CasesIndex - c(0, CasesIndex[-length(CasesIndex)]) - 1
517 R> p <- sum(Aeso[,2]) / nrow(Aeso)
518 R> p
519 [1] 0.01597444
```

520 *Scan statistic approach.* Considering the sequence of the date (`Aeso[, 2]` vector in the previous  
 521 line code), Glaz et al., 2009 give the scan statistic values for different choices of the window  
 522 length  $k$ ; the corresponding statistical significance; and the position of the window that realizes  
 523 the maximal value. They also present the method of Nagarwalla, 1996 using a scan statistic  
 524 with a variable window size for which the statistical significance is established by Monte Carlo  
 525 method. The results are presented in Table 2.

$k$	value	$p$ -value	begin	end
100	7	0.08833	1233	1305
200	10	0.04993	1233	1390
300	15	0.00141	1233	1491
365	16	0.00271	1233	1583
Nagarwalla	15	0.00580	1233	1491

**Table 2** – Results of the scan statistic approaches.

526 We can observe that the different statistical significances are very different depending on  
 527 the window size choice: Using nominal level equal to 1%, we get not significant  $p$ -values for a  
 528 window size  $k = 100$  or  $k = 200$  and significant ones for  $k = 300, 365$  and for the Nagarwall  
 529 method.

*Log Likelihood Ratio test and local score approach.* We propose here to consider an eventual drift  
 in the parameter  $p$ . Let us consider  $H_0: p = p_0$  and  $H_1: p = p_0 \cdot (1 + \delta)$  for a given  $\delta$  value. Let  
 us consider first  $\delta = 5\%$ . We associate to the Time Between Events sequence (called `tbe` in the  
 previous line code) the following score sequence computed using

$$X(tbe) = \lfloor E \cdot \ln \left( \frac{f_1(tbe)}{f_0(tbe)} \right) \rfloor$$

530 with  $f_i$  the probability for a random variable distributed as a geometrical law of parameter  $p_i$ , for  
 531  $i = 0, 1$ ;  $E$  a tuning parameter we have previously presented in Subsection .

```
532 R> p0 <- p
533 R> delta <- 0.05
534 R> p1 <- p0 * (1 + delta)
535 R> # Choice for the value of E
536 R> # in order to allow at least 3 non negative scores
537 R> E <- 1
538 R> # Maximum of the scores in a geometrical model
539 R> maxXk <- floor(E * log(p1 / p0))
540 R> while (maxXk < 3) {
541 +   E <- E+1
542 +   maxXk <- floor(E * log(p1 / p0)) }
543 R> E
544 [1] 62
```

545 This leads to the following score sequence

```
546 R> ScoreSeq <- floor(E*log(dgeom(tbe, prob = p1) /
547 +   dgeom(tbe, prob = p0)))
548 R> ScoreSeq
549 R> minX <- min(ScoreSeq)
550 R> maxX <- max(ScoreSeq)
551 [1] -6 -5 -4 1 -21 -2 -2 -3 2 3 2 2 2 1 -1
552 [16] 2 2 0 2 2 2 1 -2 -3 2 -2 -4 0 -1 2
553 [31] 2 2 1 0 1
```

554 *Score distribution.* Let us compute the score distribution by two different ways: The first one  
 555 based on the estimation on the score appearance on the observed sequence and a second way  
 556 using theoretical work on geometrical model that leads to a more accurate score distribution.

```

557 R> # Estimation on the data
558 R> ProbScoreInit <- table(ScoreSeq) / sum(table(ScoreSeq))
559 R> ProbScore <- rep(0, maxX - minX + 1)
560 R> names(ProbScore) <- minX:maxX
561 R> for (i in 1:length(ProbScoreInit)) {
562 +   w <- which(names(ProbScore) == names(ProbScoreInit)[i])
563 +   ProbScore[w] <- ProbScoreInit[i] }
564 R> fonctionProbaX <- function(x,E,p0,p1) {
565 +   calcul1 <- ((x/E)-log(p1/p0)) / log((1-p1)/(1-p0))
566 +   calcul2 <- ((x+1)/E) - log(p1/p0) / log((1-p1)/(1-p0))
567 +   prob <- pgeom(prob=p0, floor(calcul1)) -
568 +     pgeom(prob=p0, floor(calcul2))
569 +   return(prob) }
570 R> ProbScoreTheo <- rep(0, maxX - minX + 1)
571 R> names(ProbScoreTheo) <- minX:maxX
572 R> for (i in 1:length(ProbScoreTheo)) {
573 +   ProbScoreTheo[i]<-fonctionProbaX(
574 +     x = as.numeric(names(ProbScoreTheo)[i]),
575 +     E, p0, p1 = (1 + delta) * p0) }
576 R> head(cbind(ProbScore, ProbScoreTheo))
577     ProbScore  ProbScoreTheo
578 -21 0.02857143    0.0001725077
579 -20 0.00000000    0.0002380570
580 -19 0.00000000    0.0003285136
581 -18 0.00000000    0.0004533418
582 -17 0.00000000    0.0006256021
583 -16 0.00000000    0.0008132325
584 R> tail(cbind(ProbScore, ProbScoreTheo))
585     ProbScore  ProbScoreTheo
586 -2 0.11428571    0.07592718
587 -1 0.05714286    0.09869924
588  0 0.08571429    0.14228152
589  1 0.14285714    0.19634549
590  2 0.37142857    0.27095262
591  3 0.02857143    0.01597444

```

592 Let us then compute, for the given shift  $\delta$ , the local score value and its  $p$ -value with the two  
 593 different score distributions. As the length sequence is very short,  $n = 35$ , we use the exact  
 594 method to establish the  $p$ -value with the function `daudin`. The begin and end index are also  
 595 given.

```

596 R> localScoreC(ScoreSeq)

```

```

597 $localScore
598 value begin end
599 22 9 22
600 $suboptimalSegmentScores
601 value begin end
602 [1,] 1 4 4
603 [2,] 22 9 22
604 $RecordTime
605 [1] 4 9
606 R> LS <- localScoreC(ScoreSeq)$localScore[1]
607 R> BeginTbe <- localScoreC(ScoreSeq)$localScore[2]
608 R> EndTbe <- localScoreC(ScoreSeq)$localScore[3]
609 R> BeginDate <- CasesIndex[BeginTbe - 1]
610 R> EndDate <- CasesIndex[EndTbe]
611 R> pvalue <- daudin(
612 + localScore = LS,
613 + sequence_length = length(tbe),
614 + score_probabilities = ProbScoreTheo,
615 + sequence_min = minX,
616 + sequence_max = maxX )
617 R> pvalue
618 [1] 0.02647577
619 R> BeginDate
620 [1] 1233
621 R> EndDate
622 [1] 1491

```

623 The segment that realizes the local score value begins at the date index 1233 and ends at the  
624 date index 1491 which corresponds to the segment highlighted by the scan statistic approach  
625 with a window size  $k = 300$ . Its statistical significance of the observed local score is around  
626 0.026.

627 One could say that the choice of the window length in the scan statistic method does not  
628 have to be done in the local score one. But the choice is change in choosing a  $\delta$  value to construct  
629 the score function based on  $p_0$  and  $p_1$ . When no previous knowledge exists on the length of the  
630 segment we want to highlight, it is easier to choose the smallest drift we would like to detect.  
631 Let us have a look for a set of different  $\delta$  values from 1% to 5%.

632 We can observed in Table 3 that the local score value does not change and neither the seg-  
633 ment that realizes the local score: See *b.tbe* (respectively *b.date*) the begin index in the tbe (resp.  
634 date sequence and see *e.tbe* (respectively *e.date*) the end index in the tbe (resp. date) sequence.  
635 Moreover, its statistical significance is quite constant and around 3%.

636 *Direct analysis on the 0-1 sequence.* We propose below to analysis the initial sequence of occur-  
637 rences (0-1) without constructing the TBE sequence. The model is then based on a Bernoulli  
638 distribution with still parameter  $p_0 = 0.01597444$ . Let us consider a drift  $p_1 = 1.05 \cdot p_0$ . We have

$\delta$	$E$	Local score	$p$ – value	b.tbe	e.tbe	b.date	e.date
0.01	302	22	0.02962169	9	22	1233	1491
0.02	152	22	0.02835914	9	22	1233	1491
0.03	102	22	0.02757781	9	22	1233	1491
0.04	77	22	0.02735163	9	22	1233	1491
0.05	62	22	0.02647577	9	22	1233	1491

**Table 3** – Local score value, its statistical significance, the begin and end indices position in the *tbe* sequence and in the *date* sequence, for different value of  $\delta$ . We also give in the second column the tuning parameter  $E$  used to get at least three non negative scores.

639 two different score values:  $(\ln(p_1(1 - p_0)/(p_0(1 - p_1))) + \log((1 - p_1)/(1 - p_0)))$  corresponding  
640 to 1 and  $\ln((1 - p_1)/(1 - p_0))$  to 0.

```
641 R> occur <- Aeso[, 2]
```

```
642 R> p1 <- 1.05 * p0
```

```
643 R> p1
```

```
644 [1] 0.01677316
```

```
645 R> (log(p1*(1-p0) / (p0*(1-p1))) + log((1-p1) / (1-p0)))
```

```
646 [1] 0.04879016
```

```
647 R> log((1-p1) / (1-p0))
```

```
648 [1] -0.0008120179
```

649 These first score values lead us to choose a tuning parameter  $E$  equal to 1000 in order to  
650 keep the proportion between those two values and to get integer values which allow to use  
651 the exact method. We recall that this change in the score function has only consequences on  
652 the local score value, but no ones on the segment that realized the local score and neither the  
653 statistical significance.

```
654 R> E <- 1000
```

```
655 R> ScoreSeq2 <- floor(E*(occur*log(p1*(1-p0) / (p0*(1-p1)))
```

```
656 +   + log((1-p1) / (1-p0))))
```

```
657 R> table(ScoreSeq2)
```

```
658 R> localScoreC(ScoreSeq2)$localScore
```

```
659 score.seq2
```

```
660   -1   48
```

```
661 2156   35
```

```
662 $localScore
```

```
663 value begin   end
```

```
664   476 1233 1491
```

665 The two possible scores are then: -1 and 48. We observe that the segment that realizes the  
666 local score is still the same than with the geometrical model with a begin index 1233 and an end  
667 index 1491.

668 Let us have a look on the statistical significance. But first let us give the score distribution.  
669 The probabilities under  $H_0$  associated to the scores are equal to  $1 - p_0$  for -1 and  $p_0$  for 48.

```
670 R> ProbScores <- rep(0, 50)
```

```
671 R> names(ProbScores) <- -1:48
```

```

672 R> ProbScores[1] <- 1-p0
673 R> ProbScores[length(ProbScores)] <- p0
674 R> # expected score
675 R> -1*(1-p0) + 48*p0
676 [1] -0.2172524
677 R> LS <- localScoreC(ScoreSeq2)$localScore[1]
678 R> daudin(
679 +   localScore = LS,
680 +   sequence_length = length(ScoreSeq2),
681 +   score_probabilities = ProbScores,
682 +   sequence_min = -1,
683 +   sequence_max = 48 )
684 [1] 0.02497491

```

685 The  $p$ -value is of the same size as the previous study with the local score approach and with  
686 the geometrical model on the TBE sequence.

687 The two studies based on the local score highlight the same segment than the scan statistic  
688 with a window size choice of 300 and than the one highlighted by the method of Nagarwalla  
689 with a variable window. This segment is statistically significant in each method. The local score  
690 both avoids the window length choice and allows the statistical significance to be theoretically  
691 established.

## 692 Genomic regions associated with phenotypic differentiation of European local pig breeds

693 The original dataset is based on European local pig breeds genetically characterized using  
694 DNA-pool sequencing data and phenotypically characterized using breed level phenotypes re-  
695 lated to stature, fatness, growth and reproductive performance traits. It is composed of 19 popu-  
696 lations of European local pig breeds and 7 populations of industrial breeds. The genetic diversity  
697 is assessed through a SNP (Single Nucleotide Polymorphism) of medium density array leading  
698 to 16,403,270 SNPs covering 18 chromosomes of the pig genomes after filtering out SNPs with  
699 missing data. The second part of the original dataset consists of phenotype characterizations of  
700 each breed, combined into four distinct groups summarizing stature, fatness, growth, and repro-  
701 ductive performance. The purpose of the study published in Poklukar et al., 2023 is to detect  
702 genomic regions with selection signatures linked to phenotypic traits in order to uncover poten-  
703 tial candidate genes that may be under adaptation to specific environments. The methododology  
704 in Poklukar et al., 2023 uses the same approach of Coop et al., 2010 leading to elaborate a Bayes  
705 Factor measuring the link between phenotypic and genotypic variations for each SNP. Statistical  
706 significance is then assessed on a SNP by SNP base correcting the multitest problem with a False  
707 Discovery Rate (FDR) approach from Benjamini and Hochberg, 1995. They finally revealed 234  
708 regions associated with stature, fatness, growth or reproduction traits.

709 Here we propose to use a local score approach to analyze the final dataset containing the  
710 Bayes Factors associated to stature traits kindly provided by the authors of Poklukar et al., 2023.  
711 For each of the about 16 millions SNPs covering 18 pig chromosomes, we have the SNP posi-  
712 tions and the associated Bayes Factor statistics. Table 4 shows the number of points for each  
713 chromosome. Note that a Bayes Factor is a real number, and  $p$ -values associated to the local

714 score can not be directly assessed by the function `karlin`, `mcc` nor `daudin` as there associ-  
 715 ated methodologies request integer scores. A proposed solution is to discretize the scores. In  
 716 the second part of this illustration, we also assess the effect of this discretisation on the results,  
 717 comparing three schemes : 1. real scores 2. scores multiplied by 10 and rounded to closest unit  
 718 3. scores rounded to the closest unit. Due to the length of the sequence, we proceed to the eval-  
 719 uation of  $p$ -value through `karlinMonteCarlo` and `karlinMonteCarlo_double` functions  
 720 see Formula (4).

**Table 4** - Number of SNPs by pig chromosome in the dataset.

Chromosome	SNPs count
1	1427539
2	1072176
3	946332
4	923731
5	792366
6	1206701
7	899034
8	1110422
9	1020531
10	695091
11	664610
12	563400
13	1225446
14	1029162
15	899320
16	694637
17	570599
18	417965

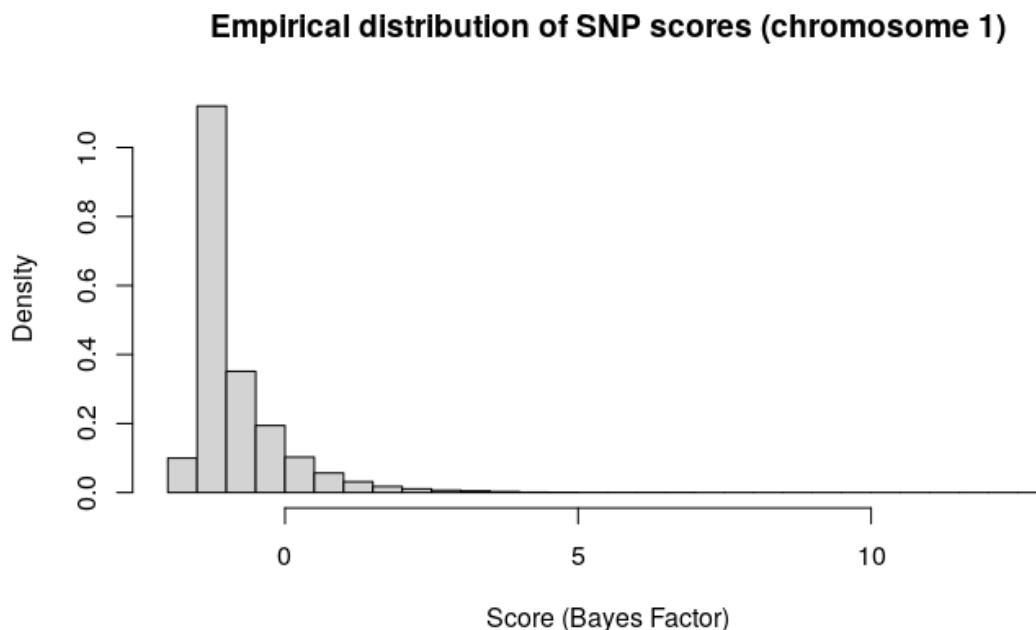
721 *Data analysis.* In order to analyze this big data file, we proceed chromosome by chromosome  
 722 and use the R library `sqldf` to load the data. Below is the R code to create a sqlite database file  
 723 'morpho\_new\_all\_Dim1.sqlite' with the data file 'morpho\_new\_all\_Dim1.flkadapt', then to load  
 724 the data of chromosome 1 (refseq id : 'NC\_010443.5') into a dataframe call 'NC\_010443.5'.

```
725 R> library(sqldf)
726 R> read.csv.sql(file="./morpho_new_all_Dim1.flkadapt",
727 +   sql = c("attach 'morpho_new_all_Dim1.sqlite' as new",
728 +   "create table new.morpho as select * from file"), sep=" ")
729 R> NC_010443.5 <- sqldf(paste(
730 +   "select * from morpho where chr='NC_010443.5'",
731 +   "AND [converge.null]='1' AND converge='1' ;",
732 +   sep=" "),
733 +   dbname = "morpho_new_all_Dim1.sqlite")
```

734 Without lost of generality, we present here the detailed analysis and results on the chromo-  
 735 some 1.

```
736 R> summary(NC_010443.5$bf)
737   Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
738 -1.9372 -1.4153 -1.1766 -0.8759 -0.6375 12.5007
```





**Figure 5** – Empirical SNP score distribution observed on the chromosome 1.

739 Note that the empirical expectation of the score is strictly negative and there is strictly posi-  
 740 tive scores as expected for a meaningful local score analysis. Figure 5 shows the empirical score  
 741 distribution obtain on chromosome 1.

742 Figure 6 shows the SNPs score observed along the chromosome 1 and the associated Lindley  
 743 process calculated as follows.

```
744 R> ScoresLindley <- lindley(NC_010443.5$bf)
```

745 The local score, its position, and all the suboptimal scores are calculated by the `localScoreC`  
 746 function.

```
747 # Calculate localscore and suboptimal localscores.
```

```
748 R> scores <- localScoreC_double(NC_010443.5$bf)
```

```
749 R> print(scores$localScore)
```

```
750     value      begin      end
751  702.7715 642973.0000 645759.0000
```

```
752 # Position of the segment realizing the local score
```

```
753 R> print(NC_010443.5$pos[scores$localScore[c("begin", "end")]])
```

```
754 [1] 86184149 86566846
```

```
755 # Length (in base-pairs) of this segment
```

```
756 R> diff(NC_010443.5$pos[scores$localScore[c("begin", "end")]])
```

```
757 [1] 382697
```

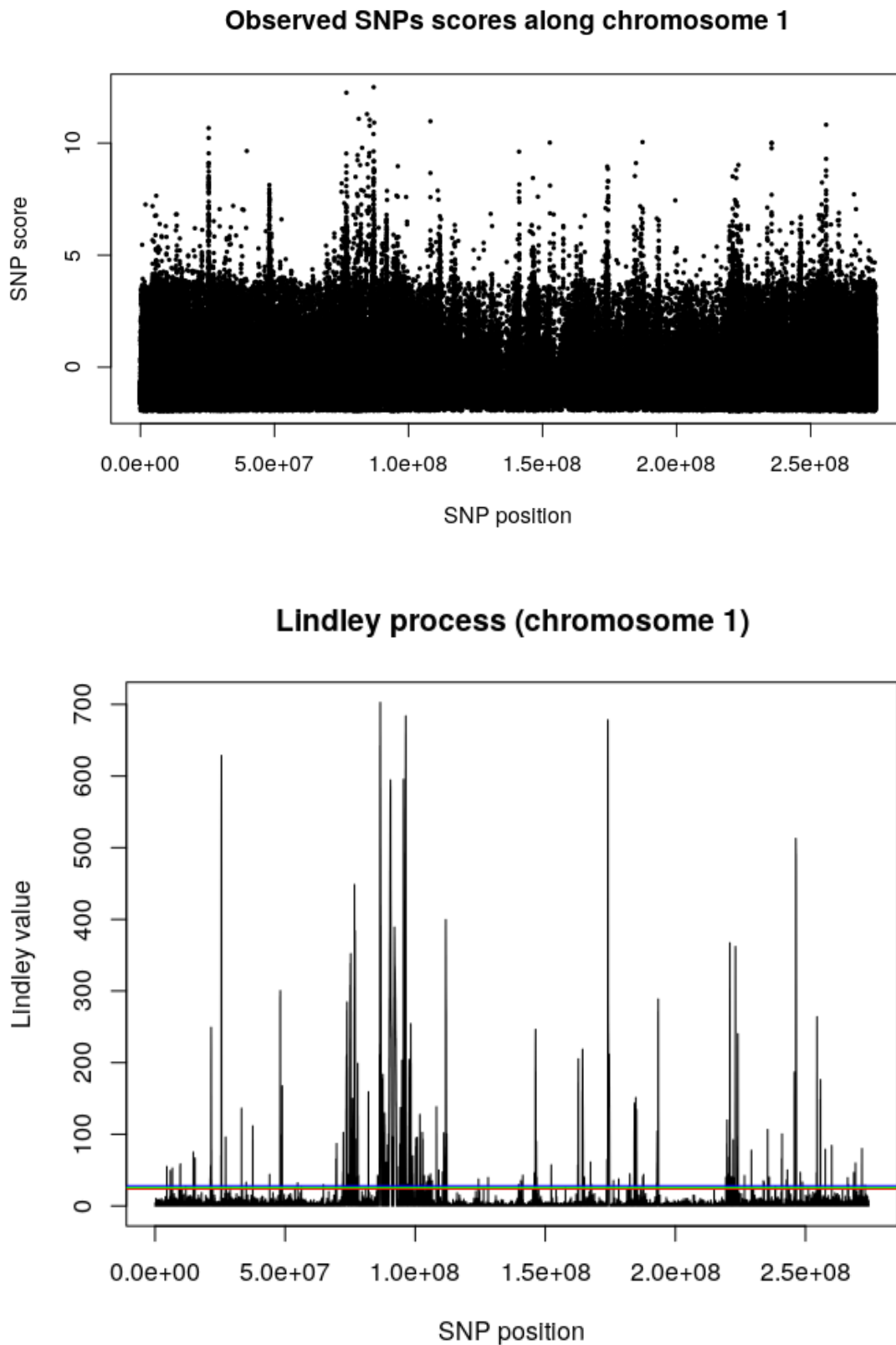
```
758 # Calculate p-value of the local score and estimation
```

```
759 # of the Karlin parameters  $K^*$  and  $\lambda$ 
```

```
760 R> ResKarlinMC <- karlinMonteCarlo_double(
```

```
761 +   local_score = scores$localScore["value"],
```

```
762 +   sequence_length = length(NC_010443.5$bf),
```



**Figure 6** – SNP scores observed along the chromosome 1: Top) SNP Score values; Bottom) Associated Lindley process, with horizontal lines representing the thresholds associated to the statistical significance of the local score at 5% (red), 1% (green) and 1‰(blue) levels.

```

763 + FUN = function(x, simulated_sequence_length)
764 +   {return(sample(x = x,
765 +     size = simulated_sequence_length,
766 +     replace = TRUE))},
767 + x = NC_010443.5$bf,
768 + simulated_sequence_length = 10000,
769 + numSim = 1000, plot = FALSE)
770 R> kStar <- ResKarlinMC$`K*`
771 R> lambda <- ResKarlinMC$lambda
772 R> print(ResKarlinMC)
773 $`p-value`
774 value
775     0
776
777 $`K*`
778 [1] 11.18694
779
780 $lambda
781 [1] 0.8285248

```

782 The local score on the chromosome 1 is 702.7715 and is realized by the segment situated in  
783 position (86184149, 86566846) with a  $p$ -value  $< 10^{-16}$ .

784 In the same way, we assess the statistical significance of the sub-optimal segments scores.  
785 As mentioned in Fariello et al., 2017, the local score threshold given for a first order risk  $\alpha$  also  
786 ensures a first order risk  $\alpha$  for at least one false positive among all excursions above this threshold.  
787 In other word all excursions above this threshold can be considered as significant sub-optimal  
788 segments scores. On the chromosome 1, we found a total of 67535 segments with positive  
789 cumulative score, from which 225 segments appear to be significant at 5%-level, 210 segments  
790 at 1%-level, and 183 segments at 1‰-level. See code below to compute these results.

```

791 R> SegmentScores <- as.data.frame(scores$suboptimalSegmentScores)
792 R> print(dim(SegmentScores)[1])
793 [1] 67535
794 R> # Sorting sub-optimal segments scores in decreasing
795 R> # order of scores
796 R> DecScoreOrder <- order(SegmentScores$value, decreasing = TRUE)
797 R> SegmentScores = SegmentScores[DecScoreOrder,]
798 R> # Calculating "p-values" majorant of each suboptimal
799 R> # segment scores until a threshold of 5% is reached.
800 R> k_star <- res.karlinMC$`K*`
801 R> lambda <- res.karlinMC$lambda
802 R> alphaMax <- 0.05
803 R> pv <- 0.0
804 R> SegmentScores$pvkarlinMC <- NA
805 R> i <- 1
806 R> while ((pv <= alphaMax) && (i <= dim(SegmentScores)[1])) {

```

```

807 + # Karlin: for n great,
808 + #P( ln(n)/lambda+x>= M) = exp(-K_star*exp(-lambda*x))
809 + # thus we set ln(n)/lambda+x = local_score and obtain
810 + # x = local_score - ln(n)/lambda
811 + x <- SegmentScores$value[i] -
812 +     log(length(NC_010443.5$bf)) / lambda
813 + # now we calculate p-value with our approximate
814 + # K star and lambda
815 + pv <- 1 - exp(-k_star * exp(-lambda * x))
816 + if (pv <= alphaMax) {
817 +     SegmentScores$pvkarlinMC[i] <- pv
818 + }
819 + i <- i + 1
820 + }
821 R> SegmentScoresSignif <-
822 +     SegmentScores[!is.na(SegmentScores$pvkarlinMC),]
823 R> # Number of significative sub-optimal segments scores at
824 R> # tthreshold 0.05, 0.01 and 0.001
825 R> alpha <- c(0.05, 0.01, 0.001)
826 R> print(sapply(alpha,
827 +     function(x){return(sum(
828 +         SegmentScoresSignif$pvkarlinMC <= x))}))
829 [1] 225 210 183

```

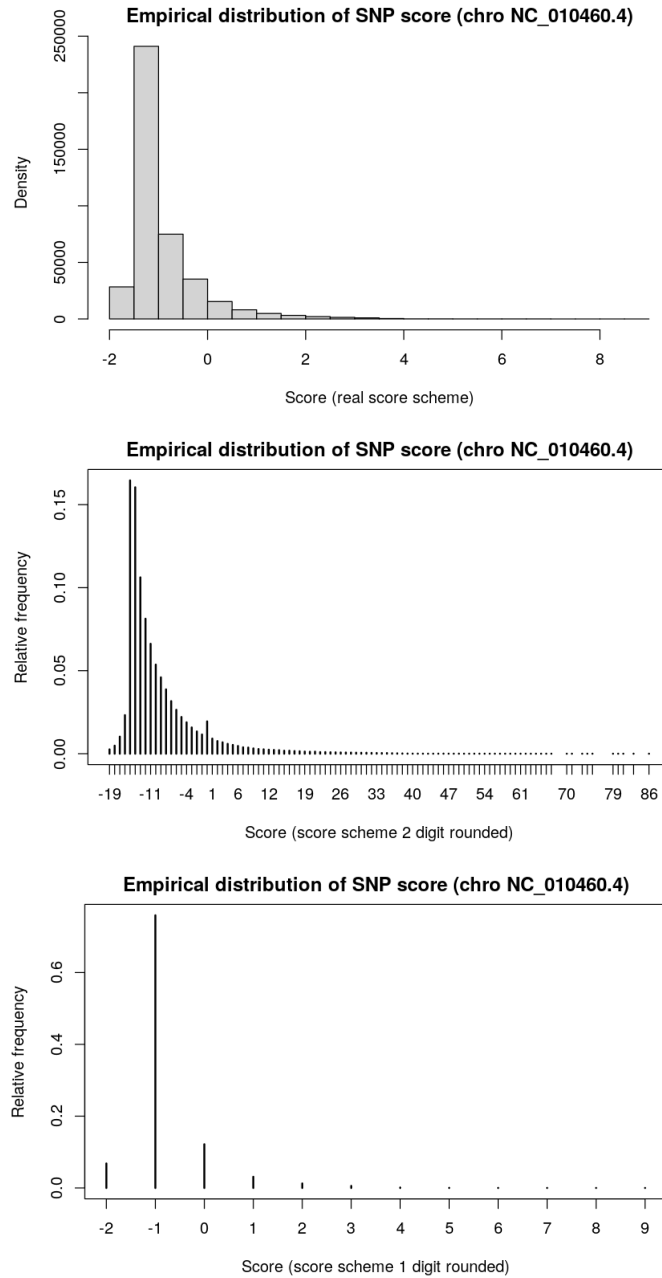
830 The chromosome 1 contains a total of 1421525 SNPs with individual scores. The proportion  
831 of SNPs presents in significative segments compared to this total, approximates 4% (see Table 5).  
832 In the paper Poklukar et al., 2023, the whole genome analysis based on the Bayes Factor retains  
833 2 segments significative at 5% corrected for multiple test, compared here to 225 segments sig-  
834 nificatively detected by the local score approach for the stature trait.

**Table 5** – Numbers and proportions of SNPs present in a significative segment according to the test threshold.

	Test threshold	0.05	0.01	0.001
Number of SNPs present in significant segment		54129	52954	51357
Proportion of SNP present in significant segment		0.04	0.04	0.04

835 *Score discretization assessment.* Three scoring schemes are compared: a) real scores as given by  
836 the Bayes Factor in input b) one decimal scores times 10 then rounded c) rounded scores to  
837 closest unit. b) and c) give integer scores. Figure 7 shows the empirical distributions of each score  
838 scheme obtained from chromosome 18. Other chromosomes show very close distributions (not  
839 shown).

840 For each chromosome of the whole genome, Table 6 summaries the number of significant  
841 detected segments applying a level of 5%, 1% and 1%. These numbers are also shown regarding  
842 the scoring scheme.



**Figure 7** – Empirical distributions of SNP scores obtained on chromosome 18 for three scoring scheme: 1. Real score 2. Two-digits rounded score 3. One-digit rounded score.

843 Let’s also look at the influence of the three scoring schemes on the length of segments that  
844 achieve the local score: Figure 8 shows comparable boxplots of the (log)-length of the detected  
845 segments below the threshold of 5% for the three scoring schemes.

846 Considering the segments obtained with the real scores as reference, Table 7 displays the  
847 numbers of false positive and false negative segments which occur with 1-digit scores and 2-  
848 digits scores. Care should be taken as a cut-off is applied to  $p$ -value less than 5% which can mod-  
849 ify the significant segments list close to the cut-off. Regarding the big range and the skewness  
850 distribution of real scores, the 2-digits rounded scores just slightly change the results, missing  
851 only 9 (0.5%) segments over 1882 real segments on the whole scale and detecting falsely 45

**Table 6** – Numbers of significant detected segments applying a level of 5%, 1% and ‰. These numbers regarding the scoring scheme are also shown.

chromosomes	Real scores			2-digits scores			Unit scores		
	5%	1%	5‰	5%	1%	5‰	5%	1%	5‰
1 NC_010443.5	226	212	185	225	213	185	218	201	179
2 NC_010444.4	70	60	49	70	59	48	67	55	44
3 NC_010445.4	56	54	46	56	53	45	56	53	45
4 NC_010446.5	90	78	70	87	78	67	89	80	71
5 NC_010447.5	94	85	67	94	85	68	105	92	72
6 NC_010448.4	93	83	74	102	91	78	109	96	88
7 NC_010449.5	94	83	71	96	85	72	96	84	72
8 NC_010450.4	131	118	108	131	118	107	127	121	112
9 NC_010451.4	132	120	105	134	123	107	145	136	120
10 NC_010452.4	108	97	86	111	98	86	123	104	91
11 NC_010453.5	67	58	50	66	57	50	60	51	44
12 NC_010454.4	52	46	40	55	47	41	56	46	41
13 NC_010455.5	161	153	130	162	156	135	160	150	133
14 NC_010456.5	135	122	105	143	126	112	140	127	107
15 NC_010457.5	142	124	110	157	138	120	152	133	113
16 NC_010458.4	74	69	64	75	71	65	81	76	69
17 NC_010459.5	70	62	57	68	61	57	68	61	57
18 NC_010460.4	87	71	61	87	73	61	86	73	63
Total	1882	1695	1478	1919	1732	1504	1938	1739	1521

852 (2.3%) segments over 1919 detected segments. Note that the brutal unit rounded score perfor-  
 853 mance essentially reflects the real segment detected with only 76 (4%) missing segments (false  
 854 negative) over 1882 and detect falsely (false positive) 157 (8%) segments over 1938 detected  
 855 segments.

**Table 7** – Considering real scoring scheme as detection reference, the table shows the numbers of segments which differ from the reference for the 2-digit rounded score scheme, and the unit rounded scheme on the whole genome analysis. "False negative": number of segment which are present in the reference, but not in considered scoring scheme; "False positive": number of segments significantly detected but not present in the reference.

score scheme	Real	2-digits	1-digit
Total detected segments (<5%)	1882	1919	1938
False negative		9 (0.5%)	76 (4%)
False positive		45 (2.3%)	157 (8%)

856

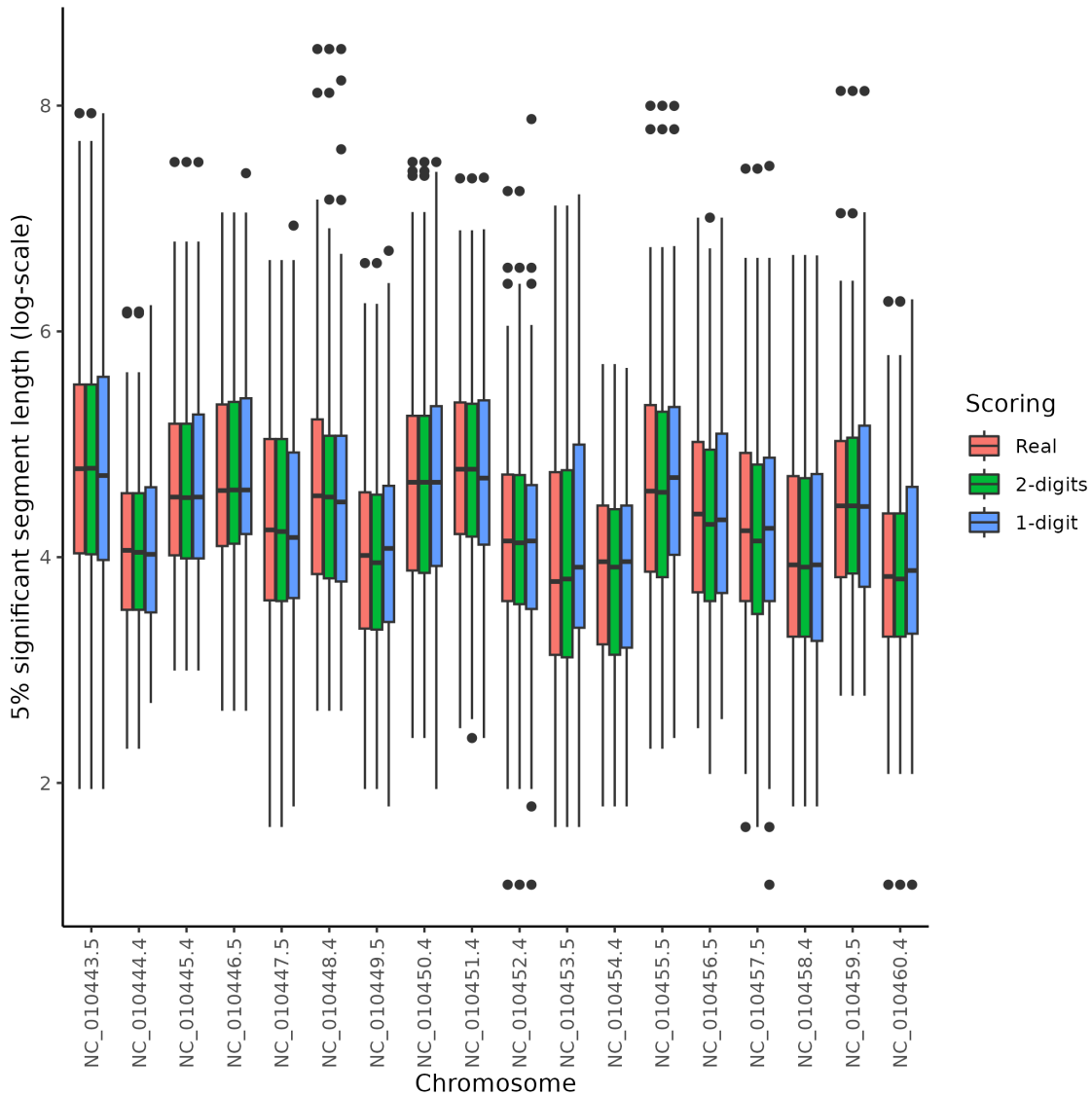
### Computational details

857 The results in this paper were obtained using R 4.3.1. R itself and all packages used are  
 858 available from the Comprehensive R Archive Network (CRAN) at [https://CRAN.R-project.](https://CRAN.R-project.org/)  
 859 [org/](https://CRAN.R-project.org/).

860

### Summary and discussion

861 When no a priori information is known about the length of the segments to be highlighted,  
 862 the local score is a dedicated tool to exploit and supplement the methods of sliding-windows or



**Figure 8** – Log-length of the detected segments below the threshold of 5% for the three scoring schemes by chromosomes.

863 scan statistics. In addition, the package allows to calculate the statistical significance and to dis-  
 864 tinguish the segments of atypical optimal scores from those appearing by chance. The package  
 865 brings together various functions that notably allow to visualize and point-out the highlighted  
 866 regions. Different ways of assessing the statistical significance are proposed. A function allows  
 867 to perform this calculation by automatically selecting the method most suited to the context  
 868 related to the length of the sequence, and the average score under a given model or learned. If  
 869 initially the local score has been defined for the identification of atypical regions within biolog-  
 870 ical sequences, it can also be useful in many fields of application as we wanted to illustrate in  
 871 our examples. It can also be applied for online analyses including breakpoint detection. Further  
 872 developments will be made for Markov models and the statistical significance of sub-optimal  
 873 segments in a future version of the package.

874

### Acknowledgements

875 The authors thank Sebastian Simon who began to build the package during his internship.

## Fundings

876

877 This work has been supported by the project “Highlight” of the *Excellence Laboratory Interna-*  
878 *tional Center of Mathematics and Computer Science in Toulouse* (Labex CIMI).

## Conflict of interest disclosure

879

880 The authors declare that they comply with the PCI rule of having no financial conflicts of  
881 interest in relation to the content of the article.

## Data, script, code, and supplementary information availability

882

883 The source code of the `localScore` package is available at: <https://cran.r-project.org/web/packages/localScore/index.html>. The package also contains the datasets analyzed in the article except the Genomic dataset that is not publicly available so far.

## References

886

887 Benjamini Y, Hochberg Y (1995). *Controlling the False Discovery Rate: A Practical and Powerful*  
888 *Approach to Multiple Testing*. *Journal of the Royal Statistical Society: Series B (Methodological)*  
889 **57**, 289–300. <https://doi.org/10.1111/j.2517-6161.1995.tb02031.x>.

890 Cellier D, Charlot F, Mercier S (2003). *An improved approximation for assessing the statistical sig-*  
891 *nificance of molecular sequence features*. *Jour. Appl. Prob.* **40**, 427–441. [https://doi.org/](https://doi.org/10.1239/jap/1053003554)  
892 [10.1239/jap/1053003554](https://doi.org/10.1239/jap/1053003554).

893 Chabriac C, Lagnoux A, Mercier S, Vallois P (2014). *Elements related to the largest complete ex-*  
894 *cursion of a reflected BM stopped at a fixed time. Application to local score*. *Stochastic Processes*  
895 *and their Applications* **124**, 4202–4223. [https://doi.org/10.1016/j.spa.2014.07.](https://doi.org/10.1016/j.spa.2014.07.003)  
896 [003](https://doi.org/10.1016/j.spa.2014.07.003).

897 Coop G, Witonsky D, Di Rienzo A, Pritchard JK (2010). *Using environmental correlations to identify*  
898 *loci underlying local adaptation*. *Genetics* **185**, 1411–1423. [https://doi.org/10.1534/](https://doi.org/10.1534/genetics.110.114819)  
899 [genetics.110.114819](https://doi.org/10.1534/genetics.110.114819).

900 Fariello M, Boitard S, Mercier S, Robelin D, Faraut T, Arnould C, Le Bihan-Duval E, Recoquilly  
901 J, Salin G, Dahais G, Pitel F, Leterrier G, Sancristobal M (2017). *Accounting for Linkage Dise-*  
902 *quilibrium in genome scans for selection without individual genotypes : the local score approach*.  
903 *Molecular Ecology* **26(14)**, 3700–3714. <https://doi.org/10.1111/mec.14141>.

904 Glaz J, Pozdnyakov V, Wallenstein S (2009). *Scan statistics - Methods and applications*. Birkhauser  
905 Bosten. Chap. 1. <https://doi.org/10.1007/978-0-8176-4749-0>.

906 Grusea S, Mercier S (2020). *Improvement on the distribution of maximal segmental score in a Mar-*  
907 *kovian sequence*. *Journal of Applied Probability* **57.1**, 29–52. [https://doi.org/10.1017/](https://doi.org/10.1017/jpr.2019.75)  
908 [jpr.2019.75](https://doi.org/10.1017/jpr.2019.75).

909 Hassenforder C, Mercier S (2007). *Exact Distribution of the Local Score for Markovian Sequences*.  
910 *AIMS* **59**, 741–755. <https://doi.org/10.1007/s10463-006-0064-6>.

911 Karlin S, Altschul SF (1990). *Methods for assessing the statistical significance of molecular sequence*  
912 *features by using general scoring schemes*. *PNAS* **87**, 2264–2268. [https://doi.org/10.](https://doi.org/10.1073/pnas.87.6.2264)  
913 [1073/pnas.87.6.2264](https://doi.org/10.1073/pnas.87.6.2264).

914 Karlin S, Dembo A (1992). *Limit distributions of maximal segmental score among Markov-dependent*  
915 *partial sums*. *AdAP* **24**, 113–140. <https://doi.org/doi.org/10.2307/1427732>.



- 916 Knox G (1959). *Secular pattern of congenital oesophageal atresia*. *British Journal of Preventive Social*  
917 *Medicine* **13**, 222–226. <https://doi.org/10.1136/jech.51.2.110>.
- 918 Kyte J, Doolittle R (1982). *A simple method for displaying the hydropathic character of a protein*.  
919 *Journal of molecular biology* **157**, 105–132. [https://doi.org/10.1016/0022-2836\(82\)](https://doi.org/10.1016/0022-2836(82)90515-0)  
920 [90515-0](https://doi.org/10.1016/0022-2836(82)90515-0).
- 921 Lagnoux A, Mercier S, Vallois P (2017). *Statistical significance based on length and position of the*  
922 *local score in a model of i.i.d. sequences*. *Bioinformatics* **33**, 654–660. [https://doi.org/10.](https://doi.org/10.1093/bioinformatics/btw699)  
923 [1093/bioinformatics/btw699](https://doi.org/10.1093/bioinformatics/btw699).
- 924 Mercier S (2020). *Transferring biological sequence analysis tools to break-point detection for on-line*  
925 *monitoring: A control chart based on the Local Score*. *Qual. Reliab. Engng. Int.* **36**, 2379–2397.  
926 <https://doi.org/doi.org/10.1002/qre.2703>.
- 927 Mercier S, Daudin J (2001). *Exact Distribution for the Local Score of One i.i.d. Random Sequence*.  
928 *Jour. Comp. Biol* **8**, 373–380. <https://doi.org/10.1089/106652701752236197>.
- 929 Nagarwalla N (1996). *A scan statistic with a variable window*. *Statistics in Medecine* **15**, 845–850.  
930 [https://doi.org/10.1002/\(sici\)1097-0258\(19960415\)15:7/9<845::aid-](https://doi.org/10.1002/(sici)1097-0258(19960415)15:7/9<845::aid-sim254>3.0.co;2-x)  
931 [sim254>3.0.co;2-x](https://doi.org/10.1002/(sici)1097-0258(19960415)15:7/9<845::aid-sim254>3.0.co;2-x).
- 932 Poklukar K, Mestre C, Škrlep M, Čandek-Potokar M, Ovilo C, Fontanesi L, Riquet J, Bovo S, Schi-  
933 *avo G, Ribani A, Muñoz M, Gallo M, Bozzi R, Charneca R, Quintanilla R, Kušec G, Mercat MJ,*  
934 *Zimmer C, Razmaite V, Araujo JP, et al. (2023). A meta-analysis of genetic and phenotypic diver-*  
935 *sity of European local pig breeds reveals genomic regions associated with breed differentiation for*  
936 *production traits*. *Genetics Selection Evolution* **55**. [https://doi.org/10.1186/s12711-](https://doi.org/10.1186/s12711-023-00858-3)  
937 [023-00858-3](https://doi.org/10.1186/s12711-023-00858-3).
- 938 R Core Team (2024). *R: A Language and Environment for Statistical Computing*. URL: [https://](https://www.R-project.org/)  
939 [www.R-project.org/](https://www.R-project.org/).
- 940 Simon S, Robelin D, Mercier S, Dejean S (2023). *localScore: Package for Sequence Analysis by Local*  
941 *Score*. R package version 1.0.11.
- 942 W.A. Shewhart (1931). *Economic Control of Quality of Manufactured Product*. D. Van Nostrand  
943 *Company, Inc., New York*.