

FunChIP: A Functional Data Analysis approach to cluster ChIP-Seq peaks according to their shapes

Alice Parodi alicecarla.parodi@polimi.it

Laura M. Sangalli

Piercesare Secchi

Simone Vantini

Marco J. Morelli

October 17, 2016

Contents

> `library(FunChIP)`

1 Introduction

The *FunChIP* package provides a set of methods for the *GRanges* class of the package *GenomicRanges* to cluster ChIP-Seq peaks according to their shapes, starting from a BAM file containing the aligned reads and a *GRanges* object with the corresponding enriched regions.

2 Input and Preprocessing

ChIP-Seq enriched regions are provided by the user in a *GRanges* object *GR*. The user must provide the BAM file containing the reads aligned on the positive and negative strands of the DNA. From the BAM file we can compute, for each region of the *GRanges* (let N be the total number of regions), the base-level coverage separately for positive and negative reads. These two count vectors are used to estimate the distance d_{pn} between positive and negative reads and then the total length of the fragments of the ChIP-Seq experiment d . In particular, we assume that the positive and negative counts measure the same signal, shifted by d_{pn} , as they are computed from the two ends of the sequencing fragments. The global length of the fragment is the sum between the length of the reads of the BAM file, r^1 , and the distance between the positive and negative coverage d_{pn}

$$d = d_{pn} + r.$$

The function `compute_fragments_length` computes, from the *GRanges* object and the BAM file, the estimated length of the fragments. Given a range for d_{pn} : $[d_{\min}; d_{\max}]$, the optimum distance d_{pn} is

$$d_{pn} = \operatorname{argmin}_{\delta \in [d_{\min}; d_{\max}]} \sum_{n=1}^N D(f_{n+}, f_{n-}^{\delta}),$$

where f_{n+} is the positive coverage function of the n -th region, and f_{n-}^{δ} is the negative coverage of the n th region, shifted by δ . The distance D is the square of the L^2 distance between the coverages, normalized by the width of the region. The definition of the L^2 distance is detailed in Section ??.

¹ If in the BAM file multiple length are present, r is estimated as the average length.

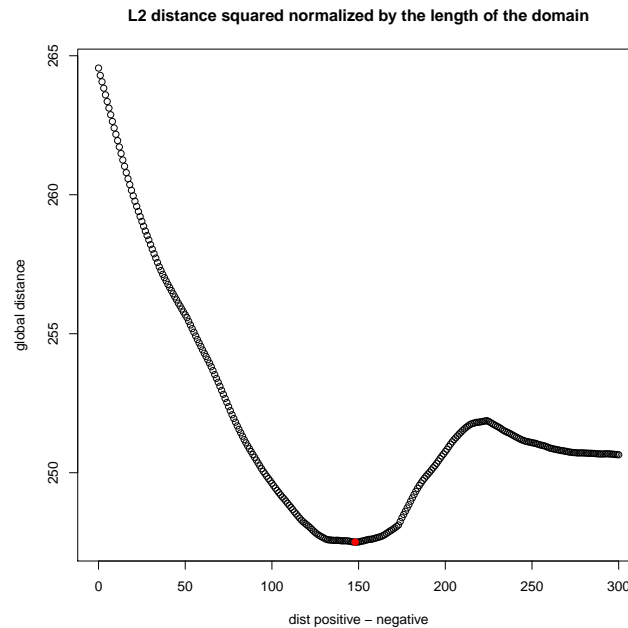


Figure 1: **Identification of d .** optimal value of d_{pn} is presented. It is the minimum of the global distance function.

```
> # load the GRanges object associated
> # to the ChIP-Seq experiment on the
> # transcription factor c-Myc in murine cells
>
> data(GR100)
> # name of the .bam file (the
> # .bam.bai index file must also be present)
>
> bamf <- system.file("extdata", "test.bam",
+                       package="FunChIP", mustWork=TRUE)
> # compute d
>
> d <- compute_fragments_length(GR, bamf, min.d = 0, max.d = 300)
[1] "estimated distance positive - negative read 148"
[1] "estimated read length 51"
> d
[1] 199
>
```

In Figure ?? the distance function is shown varying the parameter δ , and the minimum value d_{pn} is computed.

Once we have correctly identified the fragment length we can compute the final coverage function to obtain the shape of the peaks. The `pileup_peak` method for the `GRanges` class uses the BAM file to compute the base-level coverage on these regions, once the reads are extended up to their final length d . `pileup_peak` adds to the `GRanges` a counts metadata column, containing for each region a vector with length equal to the width of the region storing the coverage function.

```
> # associate to each peak
> # of the GRanges object the correspondent
```

```
> # coverage function
>
> peaks <- pileup_peak(GR, bamf, d = d)
> peaks
```

GRanges object with 100 ranges and 1 metadata column:

| | seqnames | ranges | strand | counts |
|-------|----------|----------------------|--------|--------------|
| | <Rle> | <IRanges> | <Rle> | <list> |
| [1] | chr18 | [3337524, 3338025] | * | 7,8,8,... |
| [2] | chr18 | [4369126, 4369352] | * | 7,9,9,... |
| [3] | chr18 | [4375448, 4375883] | * | 8,8,8,... |
| [4] | chr18 | [4715744, 4716162] | * | 5,5,5,... |
| [5] | chr18 | [4716374, 4716597] | * | 15,15,15,... |
| ... | ... | ... | ... | ... |
| [96] | chr18 | [35112325, 35112593] | * | 12,12,12,... |
| [97] | chr18 | [35113538, 35114826] | * | 9,8,8,... |
| [98] | chr18 | [35118063, 35118304] | * | 10,10,10,... |
| [99] | chr18 | [35182164, 35182425] | * | 8,8,8,... |
| [100] | chr18 | [35205390, 35205649] | * | 14,14,14,... |

```
seqinfo: 20 sequences from an unspecified genome; no seqlengths
>
```

Additional information can be found in the help page of the `pileup_peak` method.

3 Smoothing

The counts metadata is approximated by a combination of splines to guarantee the smoothness and regularity needed for further analysis, as described in the following Sections.

The preprocessing steps carried out in the `smooth_peak` method are the following:

- *Removal of the background and extension.* In ChIP-Seq experiments, peaks may have an additive noisy background, and the removal of this background is mandatory to compare different peaks. The background is estimated as a constant value "raising" the peak and equal to the minimum value the coverage assumes. Consequently, once the background has been removed, each peak has zero as minimum value, thus allowing the peak to be indefinitely extended with zeros, if necessary. In Section ??, how this choice affects the algorithm will be discussed.
- *Smoothing.* In order to be regular enough to computed derivatives, a peak has to be transformed in a suitable functional object, as described in Section ?. The smoothing of the count vector c is performed through the projection of c on a cubic B-spline basis $\Phi = \{\phi_1, \dots, \phi_K\}$ with a penalization on the second derivative [?]. The result is a spline approximation of the data, which is continuous on the whole domain, together with its first order derivatives. Moreover, the penalization on the second derivative allows to control the global regularity of the function avoiding over-fitting and a consequent noisy spline definition. The spline approximation $s = \sum_{k=1}^K \theta_k \phi_k$ of the count vector $c = \{c_j\}$ is defined minimizing

$$S(\lambda) = \sum_{j=1}^n [c_j - s(x_j)]^2 + \lambda \int [s''(x)]^2 dx,$$

with x_j being the relative genomic coordinate the counts. The multiplying coefficient λ quantifies the penalization on the second derivative and is chosen through the Generalized Cross Validation criteria. For each peak i the GCV_i index is computed with a leave-one-out cross validation

$$GCV_i = \left(\frac{n}{n - df(\lambda)} \right) \left(\frac{SSE_i}{n - df(\lambda)} \right)$$

and then it is summed on the whole data set to obtain the global GCV . The number of degrees of freedom $df(\lambda)$ is automatically computed from the definition of the basis Φ .

The error SSE_i can be computed either on the data (SSE_i^0) or on the derivatives (SSE_i^1), to control the regularity of the function or the regularity of the derivatives, respectively:

$$SSE_i^0 = \sqrt{\sum_{j=1}^n (c_j - s(x_j))^2} \text{ or } SSE_i^1 = \sqrt{\sum_{j=1}^{n-1} (\nabla c_j - s'(x_j))^2},$$

with ∇c_j being the finite-difference approximation of the derivative of the counts vector c for the data i : $c = c(i)$, while $s'(x_i)$ is the evaluation of the first derivative $s' = s'(i)$ on the genomic coordinates. For further details on the spline definition see the `spline` function of the `fda` package.

- *Scaling of the peaks.* This optional preprocessing step makes all the curves having the same width and area. In particular all the abscissa grid are scaled to become equal to the smallest grid throughout the data, while y -values are scaled to make areas of all the curves equal to 1.

The `smooth_peak` method approximates the counts metadata by removing the background, computing the spline and potentially defining the scaled approximation. Focusing on the spline approximation, `smooth_peak` automatically chooses the optimal λ parameter according to the GCV criteria; the user can decide whether to consider the data or the derivatives to compute the SSE .

```
> # the method smooth_peak
> # removes the background and defines the spline
> # approximation from the previously computed peaks
> # with lambda estimated from the
> # GCV on derivatives. The method spans a non-uniform
> # grid for lambda from 10^-4 to 10^12.
> # ( the grid is uniform for log10(lambda) )
>
> peaks.smooth <- smooth_peak(peaks, lambda = 10^(-4:12),
+                               subsample.data = 50, GCV.derivatives = TRUE,
+                               plot.GCV = TRUE, rescale = FALSE )

[1] "iteration on lambda: 1"
[1] "iteration on lambda: 2"
[1] "iteration on lambda: 3"
[1] "iteration on lambda: 4"
[1] "iteration on lambda: 5"
[1] "iteration on lambda: 6"
[1] "iteration on lambda: 7"
[1] "iteration on lambda: 8"
[1] "iteration on lambda: 9"
[1] "iteration on lambda: 10"
[1] "iteration on lambda: 11"
[1] "iteration on lambda: 12"
[1] "iteration on lambda: 13"
[1] "iteration on lambda: 14"
[1] "iteration on lambda: 15"
[1] "iteration on lambda: 16"
[1] "iteration on lambda: 17"
>
```

In Figure ??, the plot of the GCV for both data and derivatives is shown. From this Figure we see that the optimum value of λ , which minimizes the GCV for the derivatives, is also associated to a small value of the GCV for the data thus supporting the automatic choice.

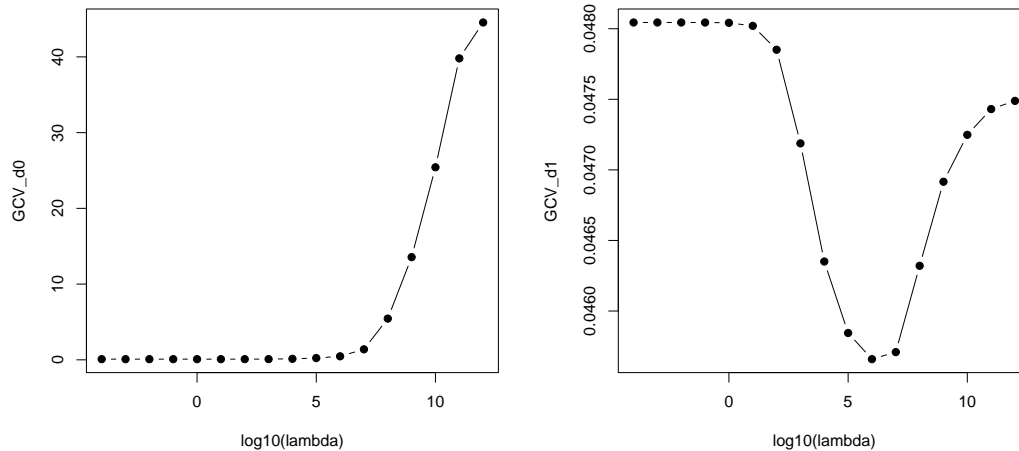


Figure 2: **Generalized Cross Validation index.** GCV computed on data (left), and on the derivatives (right), as a function of λ .

```
> # the automatic choice is lambda = 10^6
>
> peaks.smooth <- smooth_peak(peaks, lambda = 10^6,
+                             plot.GCV = FALSE)
> head(peaks.smooth)
```

GRanges object with 6 ranges and 6 metadata columns:

| | seqnames | ranges | strand | counts |
|-----|----------|--------------------|--------|--------------|
| | <Rle> | <IRanges> | <Rle> | <list> |
| [1] | chr18 | [3337524, 3338025] | * | 7,8,8,... |
| [2] | chr18 | [4369126, 4369352] | * | 7,9,9,... |
| [3] | chr18 | [4375448, 4375883] | * | 8,8,8,... |
| [4] | chr18 | [4715744, 4716162] | * | 5,5,5,... |
| [5] | chr18 | [4716374, 4716597] | * | 15,15,15,... |
| [6] | chr18 | [4921270, 4921506] | * | 8,8,8,... |

spline
<list>

| | |
|-----|---|
| [1] | 0.107576984737605,0.115943638884655,0.124571296466484,... |
| [2] | 0.106531272260595,0.125615788803367,0.145440605226015,... |
| [3] | 0.103757427817775,0.127648996987936,0.15250177579018,... |
| [4] | 0.108024633674293,0.12118157243357,0.134813515447005,... |
| [5] | 0.109476254160459,0.12659566766065,0.144359995969534,... |
| [6] | 0.104671704031439,0.120304286111549,0.136527143979231,... |

spline_der width_spline
<list> <integer>

| | | |
|-----|---|-----|
| [1] | 0.00823718792931595,0.00849663811461222,0.00875919479887256,... | 593 |
| [2] | 0.0187191316937831,0.0194522839372348,0.020199731453538,... | 335 |
| [3] | 0.0234173093333372,0.0243690129025921,0.025339543308228,... | 508 |
| [4] | 0.0129217547880229,0.0133932818084443,0.0138717632963375,... | 523 |
| [5] | 0.0168009286185096,0.0174398846432053,0.0180907582358953,... | 315 |
| [6] | 0.0153411181741832,0.0159258829799662,0.0165216697493268,... | 332 |

start_spline end_spline
<numeric> <numeric>

```
[1]      3337483      3338075
[2]      4369075      4369409
[3]      4375417      4375924
[4]      4715698      4716220
[5]      4716303      4716617
[6]      4921234      4921565
```

```
-----
```

```
seqinfo: 20 sequences from an unspecified genome; no seqlengths
```

```
> # maintaining this choice of lambda smooth_peak
> # can also define the scaled approximation
> # of the spline
>
> peaks.smooth.scaled <- smooth_peak(peaks, lambda = 10^6,
+                                   plot.GCV = FALSE, rescale = TRUE)
> head(peaks.smooth.scaled)
```

```
GRanges object with 6 ranges and 8 metadata columns:
```

| | seqnames | ranges | strand | counts |
|-----|----------|--------------------|--------|--------------|
| | <Rle> | <IRanges> | <Rle> | <list> |
| [1] | chr18 | [3337524, 3338025] | * | 7,8,8,... |
| [2] | chr18 | [4369126, 4369352] | * | 7,9,9,... |
| [3] | chr18 | [4375448, 4375883] | * | 8,8,8,... |
| [4] | chr18 | [4715744, 4716162] | * | 5,5,5,... |
| [5] | chr18 | [4716374, 4716597] | * | 15,15,15,... |
| [6] | chr18 | [4921270, 4921506] | * | 8,8,8,... |

```
spline
<list>
```

```
[1] 0.107576984737605,0.115943638884655,0.124571296466484,...
[2] 0.106531272260595,0.125615788803367,0.145440605226015,...
[3] 0.103757427817775,0.127648996987936,0.15250177579018,...
[4] 0.108024633674293,0.12118157243357,0.134813515447005,...
[5] 0.109476254160459,0.12659566766065,0.144359995969534,...
[6] 0.104671704031439,0.120304286111549,0.136527143979231,...
```

```
spline_der width_spline
<list> <integer>
```

```
[1] 0.00823718792931595,0.00849663811461222,0.00875919479887256,... 593
[2] 0.0187191316937831,0.0194522839372348,0.020199731453538,... 335
[3] 0.0234173093333372,0.0243690129025921,0.025339543308228,... 508
[4] 0.0129217547880229,0.0133932818084443,0.0138717632963375,... 523
[5] 0.0168009286185096,0.0174398846432053,0.0180907582358953,... 315
[6] 0.0153411181741832,0.0159258829799662,0.0165216697493268,... 332
```

```
start_spline end_spline
```

| | <numeric> | <numeric> |
|-----|-----------|-----------|
| [1] | 3337483 | 3338075 |
| [2] | 4369075 | 4369409 |
| [3] | 4375417 | 4375924 |
| [4] | 4715698 | 4716220 |
| [5] | 4716303 | 4716617 |
| [6] | 4921234 | 4921565 |

```
[1] 0.00263609710663517,0.0030440063908373,...
[2] 0.00122834484766426,0.00147221033932346,0.00172653795382648,0.00199155049789865,0.002267...
[3] 0.000864564256647586,0.00120333309941394,0.00156495680464632,0.00195011992527804,0.002359...
```

```
[4] 0.000788310894250635,0.000956445833827378,0.00113504914292024,0.0013243784475671,0.001524
[5] 0.00235269135248368,0.00273549871311195,0.00313329749625751,0.00354638183650247,0.00397503469483484,
[6] 0.00139356234463668,0.0016220858943517,0.00186006713373021,0.00210769852286914,0.002365173875
```

```
[1]
[2] 0.000238708662407768,0.000249059436995862,0.000259632908095442,0.00026985979570835,0.000279979570835,
[3] 0.000327571761954349,0.000350081098788865,0.000373281486886386,0.000397503469483484,0.00042190469042,
[4] 0.00044638183650247,0.00047039775,0.00049483484,0.00051935248368,0.000543668,
[5] 0.000375410625363164,0.00039025358389014,0.00040539347039775,0.00042082190469042,0.000436530506,
[6] 0.000223858745449888,0.000233220374263462,0.000242774124776864,0.000252521012181378,0.000262462051668
```

```
-----
seqinfo: 20 sequences from an unspecified genome; no seqlengths
```

```
>
```

Now the *GRanges* object contains, besides counts, 5 new metadata columns with the spline approximation evaluated on the base-level grid, its derivatives, the width of the spline and the new starting and ending points (see Figure ??). For a more detailed description of the metadata columns, see the help page of the `smooth_peak` method.

With the introduction of the smoothing, counts at the edges of the peak are connected with regularity to 0, and therefore new values different from zeros may be introduced. In order to maintain regularity, the grid is extended up to the new boundaries.

Adding to `smooth_peak` the option `rescale = TRUE` the method, beside the 5 metadata columns previously introduced, returns 2 more metadata columns with the scaled approximation of the spline and its derivatives.

Once the spline approximation is defined, the summit of the smoothed peak (or even of the scaled peak), i.e. of its spline approximation, can be detected. The summit will be used to initialize the peak alignment procedure, described in Section ??, and it can either be a user-defined parameter, stored in a vector of the same length of the GR, or automatically computed as the maximum height of the spline. The summit is stored in the new metadata column `summit_spline`. If the `rescale` option is set to `TRUE` the summit of the scaled approximation is also returned in the metadata column `summit_spline_rescaled`.

```
> # peaks.summit identifies the maximum point
> # of the smoothed peaks
>
> peaks.summit <- summit_peak(peaks.smooth)
> head(peaks.summit)
```

GRanges object with 6 ranges and 7 metadata columns:

| | seqnames | ranges | strand | counts | spline | spline_der | width_spline |
|-----|----------|--------------------|--------|--------------|---|------------|--------------|
| | <Rle> | <IRanges> | <Rle> | <list> | <list> | | |
| [1] | chr18 | [3337524, 3338025] | * | 7,8,8,... | 0.107576984737605,0.115943638884655,0.124571296466484,... | | |
| [2] | chr18 | [4369126, 4369352] | * | 7,9,9,... | 0.106531272260595,0.125615788803367,0.145440605226015,... | | |
| [3] | chr18 | [4375448, 4375883] | * | 8,8,8,... | 0.103757427817775,0.127648996987936,0.15250177579018,... | | |
| [4] | chr18 | [4715744, 4716162] | * | 5,5,5,... | 0.108024633674293,0.12118157243357,0.134813515447005,... | | |
| [5] | chr18 | [4716374, 4716597] | * | 15,15,15,... | 0.109476254160459,0.12659566766065,0.144359995969534,... | | |
| [6] | chr18 | [4921270, 4921506] | * | 8,8,8,... | 0.104671704031439,0.120304286111549,0.136527143979231,... | | |

```

                                <list>      <integer>
[1] 0.00823718792931595,0.00849663811461222,0.00875919479887256,...      593
[2] 0.0187191316937831,0.0194522839372348,0.020199731453538,...      335
[3] 0.0234173093333372,0.0243690129025921,0.025339543308228,...      508
[4] 0.0129217547880229,0.0133932818084443,0.0138717632963375,...      523
[5] 0.0168009286185096,0.0174398846432053,0.0180907582358953,...      315
[6] 0.0153411181741832,0.0159258829799662,0.0165216697493268,...      332
  start_spline end_spline summit_spline
    <numeric>  <numeric>      <integer>
[1]    3337483    3338075         444
[2]    4369075    4369409         186
[3]    4375417    4375924         174
[4]    4715698    4716220         310
[5]    4716303    4716617         121
[6]    4921234    4921565         169
-----
seqinfo: 20 sequences from an unspecified genome; no seqlengths
> # peaks.summit can identify also the maximum
> # point of the scaled approximation
>
> peaks.summit.scaled <- summit_peak(peaks.smooth.scaled,
+                                   rescale = TRUE)
> head(peaks.summit.scaled)

GRanges object with 6 ranges and 10 metadata columns:
      seqnames      ranges strand |      counts
      <Rle>        <IRanges>  <Rle> |      <list>
[1]   chr18 [3337524, 3338025]   * |    7,8,8,...
[2]   chr18 [4369126, 4369352]   * |    7,9,9,...
[3]   chr18 [4375448, 4375883]   * |    8,8,8,...
[4]   chr18 [4715744, 4716162]   * |    5,5,5,...
[5]   chr18 [4716374, 4716597]   * |   15,15,15,...
[6]   chr18 [4921270, 4921506]   * |    8,8,8,...

                                spline
                                <list>
[1] 0.107576984737605,0.115943638884655,0.124571296466484,...
[2] 0.106531272260595,0.125615788803367,0.145440605226015,...
[3] 0.103757427817775,0.127648996987936,0.15250177579018,...
[4] 0.108024633674293,0.12118157243357,0.134813515447005,...
[5] 0.109476254160459,0.12659566766065,0.144359995969534,...
[6] 0.104671704031439,0.120304286111549,0.136527143979231,...

                                spline_der width_spline
                                <list>      <integer>
[1] 0.00823718792931595,0.00849663811461222,0.00875919479887256,...      593
[2] 0.0187191316937831,0.0194522839372348,0.020199731453538,...      335
[3] 0.0234173093333372,0.0243690129025921,0.025339543308228,...      508
[4] 0.0129217547880229,0.0133932818084443,0.0138717632963375,...      523
[5] 0.0168009286185096,0.0174398846432053,0.0180907582358953,...      315
[6] 0.0153411181741832,0.0159258829799662,0.0165216697493268,...      332
  start_spline end_spline
    <numeric>  <numeric>
[1]    3337483    3338075
[2]    4369075    4369409
[3]    4375417    4375924

```



```
[4]      4715698      4716220
[5]      4716303      4716617
[6]      4921234      4921565
```

```
[1]                                     0.00263609710663517,0.0030440063908373
[2]      0.00122834484766426,0.00147221033932346,0.00172653795382648,0.00199155049789865,0.002267
[3]      0.000864564256647586,0.00120333309941394,0.00156495680464632,0.00195011992527804,0.002359
[4]      0.000788310894250635,0.000956445833827378,0.00113504914292024,0.0013243784475671,0.001524
[5] 0.00235269135248368,0.00273549871311195,0.00313329749625751,0.00354638183650247,0.00397503469483484,0
[6]      0.00139356234463668,0.0016220858943517,0.00186006713373021,0.00210769852286914,0.002365173875
```

```
[1]
[2]      0.000238708662407768,0.000249059436995862,0.000259632908095442,0.00026985979570835,0.000279
[3]      0.000327571761954349,0.000350081098788865,0.000373281486886386,0.000397571761954349,0.000419
[4]      0.000162985979570835,0.00017332651761954349,0.0001837671761954349,0.0001941081761954349,0.000204449
[5]      0.000375410625363164,0.00039025358389014,0.00040539347039775,0.00042082190469042,0.000436530506
[6] 0.000223858745449888,0.000233220374263462,0.000242774124776864,0.000252521012181378,0.000262462051668
```

```
summit_spline_rescaled summit_spline
      <integer>      <integer>
[1]      227      444
[2]      168      186
[3]      104      174
[4]      180      310
[5]      117      121
[6]      155      169
```

```
-----
seqinfo: 20 sequences from an unspecified genome; no seqlengths
```

4 The k-mean alignment algorithm and the cluster_peak method

The k-mean alignment algorithm is an efficient method to classify functional data allowing for general transformation of abscissae [?]; this general method is implemented in the package *fdakma* and various applications to real dataset are introduced in [?], [?], [?].

In particular, given

- a set of curves s_1, \dots, s_n ,
- the number of clusters K ,
- a distance function $d(s_i, s_j)$ between two curves s_i and s_j , as for example the integral of the difference $s_i - s_j$,
- a family of warping functions \mathcal{W} to transform the abscissae of the curves and therefore align the peaks. Generally, \mathcal{W} is the set of shifts or dilations or affine transformations (shift + dilation),

the algorithm, presented in Algorithm ??, is an iterative procedure to split the curves into K clusters. The introduction of the warping function $h \in \mathcal{W}$ allows each curve to be shifted, dilated, or both, to define the minimum distance between curves. The new curve $s \circ h$ has the same values of s , but its abscissa grid is modified. For example, in Figure ?? two peaks are presented: in the left panel, they are not aligned, while the right panel shows the effects of alignment; the transformation of the abscissae (shift transformation) makes the two peaks more similar, and the distance d is not anymore affected by artificial phase distance. The code generating Figure ?? calls `cluster_peak` and `plot_peak`, which are described in Section ?? and Section ??.

```
> # representation of two peaks
>
```

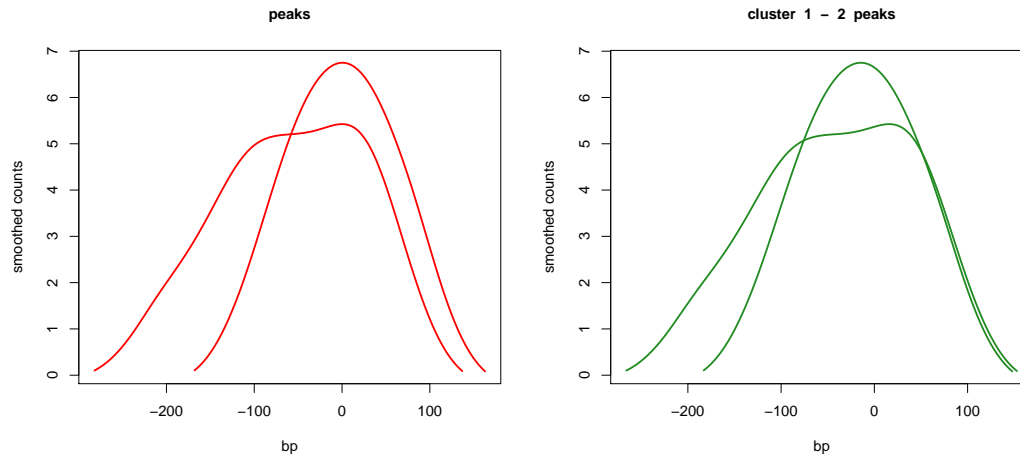


Figure 3: **Alignment procedure.** Representation of two smoothed peaks. In the left panel they are not aligned, while in the right panel they are aligned with an integer shift.

```
> par (mfrow = c(1,2))
> plot_peak(peaks.summit, index = c(6,7), col=c('red',2))
> aligned.peaks <- cluster_peak(peaks.summit[c(6,7)], parallel = FALSE ,
+                               n.clust = 1, seeds = 1, shift.peak = TRUE,
+                               weight = 1, alpha = 1, p = 2, t.max = 2,
+                               plot.graph.k = FALSE, verbose = FALSE)
> aligned.peaks
```

GRanges object with 2 ranges and 10 metadata columns:

| | seqnames | ranges | strand | counts |
|-----|----------|--------------------|--------|--------------|
| | <Rle> | <IRanges> | <Rle> | <list> |
| [1] | chr18 | [4921270, 4921506] | * | 8,8,8,... |
| [2] | chr18 | [5078473, 5078803] | * | 11,11,11,... |

```

                                spline
                                <list>
[1] 0.104671704031439,0.120304286111549,0.136527143979231,...
[2] 0.100641001405493,0.110866046987915,0.121441682487124,...

                                spline_der width_spline
                                <list>      <integer>
[1] 0.0153411181741832,0.0159258829799662,0.0165216697493268,...      332
[2] 0.0100516245524145,0.0103994035766233,0.0107528043859863,...      420
  start_spline end_spline summit_spline cluster_shift coef_shift      dist_shift
    <numeric>  <numeric>    <integer>    <list>      <list>      <list>
[1]      4921234      4921565         169          1        -15 0.447149205202752
[2]      5078418      5078837         283          1         16          0
-----
seqinfo: 20 sequences from an unspecified genome; no seqlengths
> # shift coefficients
> aligned.peaks$coef_shift

[[1]]
[1] -15

[[2]]
```

[1] 16

```
> plot_peak(aligned.peaks, col = 'forestgreen',
+           shift = TRUE, k = 1, cluster_peak = TRUE,
+           line.plot = 'spline')
```

For the specific case of ChIP-Seq data, the admitted warping functions for the k-mean alignment algorithm (in the `cluster_peak` method), are integer shifts:

$$\mathcal{W} = \{h : h(t) = t + q \text{ with } q \in \mathbb{Z}\}. \quad (1)$$

In other words, with this choice, peaks can be shifted by integer values in the *alignment* procedure of the algorithm.

Algorithm 1: k-mean alignment algorithm

Given a set of functions s_1, \dots, s_n and a number K of clusters

Template: random choice (if not provided) of the initial centers of the clusters c_1, \dots, c_k

while decrease of the distance higher than a fixed threshold **do**

foreach $i \in 1 : n$ **do**

Alignment: s_i is aligned to each template c_k : the optimal warping function $h_{i,k}^*$ in \mathcal{W} is detected

$$h_{i,k}^* = \operatorname{argmin}_{h \in \mathcal{W}} d(c_k, x_i \circ h)$$

 with the corresponding distance $d_{i,k}^* = \min_{h \in \mathcal{W}} d(c_k, x_i \circ h)$

Assignment: s_i is assigned to the best cluster

$$k_i^* = \operatorname{argmin}_{k \in 1:K} d_{i,k}^*$$

end

foreach $k \in 1 : K$ **do**

Template: identification of the new template of the cluster c_k

Normalization: the average warping function of the curves belonging to k is set to be the identity transformation

$$h(s) = s$$

end

end

In the `cluster_peak` method the distance between two curves s_1 and s_2 is defined as

$$\begin{aligned} d(s_1, s_2) &= (1 - \alpha) d_0(s_1, s_2) + \alpha w d_1(s_1, s_2) = \\ &= (1 - \alpha) \|s_1^e - s_2^e\|_p + \alpha w \|(s_1^e)' - (s_2^e)'\|_p, \end{aligned} \quad (2)$$

where

- $\|f\|_p$ is the p norm of f . In particular, for $p = 0$, $\|\cdot\|_p$ is the L^∞ norm

$$\|f\|_0 = \|f\|_{L^\infty} = \max_{x \in U} |f(x)|,$$

with U being the domain of f .

For $p = 1$, $\|\cdot\|_p$ is the L^1 norm

$$\|f\|_1 = \|f\|_{L^1} = \int_U |f(x)| dx.$$

And for $p = 2$, $\|\cdot\|_p$ is the L^2 norm

$$\|f\|_2 = \|f\|_{L^2} = \int_U (f(x))^2 dx.$$

- s_1^e and s_2^e are the functions s_1 and s_2 extended with zeros where not defined, after their backgrounds have been removed (see Section ??). The distance function is computed on the union of the domains of s_1 and s_2 (U); s_1 and s_2 need to be extended to cover the whole U .
- $\alpha \in [0, 1]$ is a coefficient tuning the contributions of the norm of the data and the norm of the derivatives. If $\alpha = 0$, the distance is computed on the data, while if $\alpha = 1$ it is based on the derivatives. Intermediate values balance these two contributions: increasing the relevance given to the derivatives emphasizes the shapes of the peaks, while data are more related to the height.
- w is a weight coefficient, essential to make the norm of the data and of the derivatives comparable. It can be user defined or computed inside the `cluster_peak` method. A suggestion for computing the weight w is given in Section ??.

4.1 Definition of weight in the distance function

If not provided, the method `cluster_peak` defines w as

$$w = \text{median} \left(\frac{d_0(s_i, s_j)}{d_1(s_i, s_j)} \right)$$

where $d_0(i, j) = \|s_i^e - s_j^e\|_p$ and $d_1(i, j) = \|(s_1^e)' - (s_2^e)'\|_p$. These matrices can be automatically computed with the `distance_peak` function.

```
> # compute the weight from the first 10 peaks
>
> dist_matrix <- distance_peak(peaks.summit)
[1] 100 1336
[1] 100 1336

> # dist matrix contains the two matrices d_0(i,j)
> # and d_1(i,j), used to compute w
> names(dist_matrix)
[1] "dist_matrix_d0" "dist_matrix_d1"

> ratio_norm <- dist_matrix$dist_matrix_d0 / dist_matrix$dist_matrix_d1
> ratio_norm_upper_tri <- ratio_norm[upper.tri(ratio_norm)]
> summary(ratio_norm_upper_tri)

  Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
 28.21  91.36  117.80  124.10  148.20  330.90

> # suggestion: use the median as weight
> w <- median(ratio_norm_upper_tri)
> w
[1] 117.7907
>
```

4.2 The cluster_peak method

The two main characteristics of the k-mean alignment algorithm used in *FunChIP* are the distance function d (defined in Equation (??)), used to compute the distance between curves, and the set of warping functions \mathcal{W} (defined in Equation (??)) considered for the alignment. The `cluster_peak` method applies the k-mean alignment algorithm with these specifications to the set of peaks stored in the *GRanges* object. In particular, the parameters `weight`², `alpha` and `p`

² `weight` can be also set to `NULL` and it will be automatically computed as specified in Section ?? . To save computational time, it is generally computed on a random sub-sample of data, whose size is set by the `subsample.weight` parameter.

define the distance used in the algorithm, while `t.max` sets the maximum shift of each peak in each iteration (in this particular case, q of Equation (??) does not vary in the whole \mathbb{Z} but $q \in \{-t.max \cdot |U|, \dots, +t.max \cdot |U|\}$, with $|U|$ being the maximum width of the spline approximation of the peaks.

Given a *GRanges* GR containing the metadata columns computed from the `smooth_peak` method, `cluster_peak` applies the k-mean alignment algorithm for all the values of k between 1 and `n.clust` (parameter of the function).

The algorithm can be run in parallel, setting to `TRUE` the `parallel` argument of the method and providing the number of cores `num.cores`. With these settings, the different applications of the algorithm, corresponding to different numbers of clusters, are executed in parallel.

As detailed in the help, the `cluster_peak` method has 2 outputs:

- The *GRanges* object, updated with new metadata columns associated to the classification. In particular, in the general case of classification with and without alignment, columns with information on the clustering of the peaks (`cluster_shift` and `cluster_NOshift`), the corresponding shifts (`coef_shift`) and the distances from the template of the clusters (`dist_shift` and `dist_NOshift`) are added.
- The graph of the global distance within clusters³ as a function of the number of clusters (if `plot.graph.k = TRUE`). This plot can be used to identify the optimal number of clusters of the partition of the data set and the effect of the alignment procedure. In particular, if `shift = NULL`, the algorithm is run both with and without alignment and two trend lines are plotted: the black line corresponds to the global distance without the shift, and the red line corresponds to the distance obtained with alignment. If `shift` is set to `TRUE` or `FALSE`, just one type of algorithm is run and the correspondent curve is plotted. For each trend line, this graph allows the identification of the optimal value of the number of clusters: for this value, the distance significantly decreases with respect to the lower values of k , and negligibly increases with respect to higher values of k (elbow in the line). The gap between the red and the black line, instead, shows the decrease of the distance when the shift is introduced.

It is relevant to point out that the algorithm can be run both on the original data and on the scaled peaks, depending on the focus of the analysis. The logic parameter `rescale` allows the user to choose.

```
> # classification of the smooth peaks in different
> # numbers of clusters, from 1 ( no distinction, only shift )
> # to 6.
>
> # here the analysis is run on the spline approximation
> # without scaling
> peaks.cluster <- cluster_peak(peaks.summit, parallel = FALSE , seeds=1:6,
+                               n.clust = 1:6, shift = NULL,
+                               weight = 1, alpha = 1, p = 2, t.max = 2,
+                               plot.graph.k = TRUE, verbose = FALSE)
> head(peaks.cluster)
```

GRanges object with 6 ranges and 12 metadata columns:

| | seqnames | ranges | strand | counts |
|-----|----------|--------------------|--------|--------------|
| | <Rle> | <IRanges> | <Rle> | <list> |
| [1] | chr18 | [3337524, 3338025] | * | 7,8,8,... |
| [2] | chr18 | [4369126, 4369352] | * | 7,9,9,... |
| [3] | chr18 | [4375448, 4375883] | * | 8,8,8,... |
| [4] | chr18 | [4715744, 4716162] | * | 5,5,5,... |
| [5] | chr18 | [4716374, 4716597] | * | 15,15,15,... |
| [6] | chr18 | [4921270, 4921506] | * | 8,8,8,... |

| | spline |
|-----|---|
| | <list> |
| [1] | 0.107576984737605,0.115943638884655,0.124571296466484,... |
| [2] | 0.106531272260595,0.125615788803367,0.145440605226015,... |
| [3] | 0.103757427817775,0.127648996987936,0.15250177579018,... |

³sum over all the peaks of the distance of each peak from the corresponding template.

```

[4] 0.108024633674293,0.12118157243357,0.134813515447005,...
[5] 0.109476254160459,0.12659566766065,0.144359995969534,...
[6] 0.104671704031439,0.120304286111549,0.136527143979231,...

                                spline_der width_spline
                                <list>    <integer>
[1] 0.00823718792931595,0.00849663811461222,0.00875919479887256,...      593
[2] 0.0187191316937831,0.0194522839372348,0.020199731453538,...      335
[3] 0.0234173093333372,0.0243690129025921,0.025339543308228,...      508
[4] 0.0129217547880229,0.0133932818084443,0.0138717632963375,...      523
[5] 0.0168009286185096,0.0174398846432053,0.0180907582358953,...      315
[6] 0.0153411181741832,0.0159258829799662,0.0165216697493268,...      332
start_spline end_spline summit_spline cluster_N0shift
<numeric>    <numeric>    <integer>    <list>
[1] 3337483    3338075      444      1,1,1,...
[2] 4369075    4369409      186      1,2,2,...
[3] 4375417    4375924      174      1,2,2,...
[4] 4715698    4716220      310      1,2,2,...
[5] 4716303    4716617      121      1,1,1,...
[6] 4921234    4921565      169      1,1,2,...
                                dist_N0shift cluster_shift
                                <list>    <list>
[1] 0.624908234181997,0.225283720971178,0.225283720971178,...      1,1,1,...
[2] 0.436709932429605,0.517291820291409,0.432369848037446,...      1,2,2,...
[3] 0.378181153745047,0.485135818958827,0.382850137813199,...      1,2,2,...
[4] 0.539505144001724,0.5547517314999,0.578858851697885,...      1,2,2,...
[5] 0.534987295842725,0.217738151115933,0.217738151115933,...      1,1,1,...
[6] 0.293729665431876,0.414829041487123,0.292960162760959,...      1,1,2,...
                                coef_shift    dist_shift
                                <list>    <list>
[1] 6,18,23,... 0.611798948745457,0.217935827035599,0.217935827035599,...
[2] 16,21,22,... 0.296361824245803,0.492957766323058,0.252045266845235,...
[3] -11,-3,-5,... 0.371425430479149,0.382138694017752,0.372516258232315,...
[4] -13,-6,-10,... 0.531584502689481,0.556766222517974,0.55605805078994,...
[5] -14,9,14,... 0.528633634301522,0.217563128269193,0.217563128269193,...
[6] 1,14,6,... 0.273367776201618,0.413050475331536,0.273737759495666,...
-----

```

```
seqinfo: 20 sequences from an unspecified genome; no seqlengths
```

```

>
> # here the analysis is run on the spline approximation
> # with scaling
> peaks.cluster.scaled <- cluster_peak(peaks.summit.scaled, parallel = FALSE , seeds=1:6,
+                                     n.clust = 1:6, shift = NULL,
+                                     weight = 1, alpha = 1, p = 2, t.max = 2,
+                                     plot.graph.k = TRUE, verbose = FALSE,
+                                     rescale = TRUE)
> head(peaks.cluster.scaled)

```

```
GRanges object with 6 ranges and 15 metadata columns:
```

```

      seqnames      ranges strand |      counts
      <Rle>        <IRanges> <Rle> |      <list>
[1] chr18 [3337524, 3338025]   * | 7,8,8,...
[2] chr18 [4369126, 4369352]   * | 7,9,9,...
[3] chr18 [4375448, 4375883]   * | 8,8,8,...
[4] chr18 [4715744, 4716162]   * | 5,5,5,...

```

```

[5] chr18 [4716374, 4716597] * | 15,15,15,...
[6] chr18 [4921270, 4921506] * | 8,8,8,...

                                spline
                                <list>
[1] 0.107576984737605,0.115943638884655,0.124571296466484,...
[2] 0.106531272260595,0.125615788803367,0.145440605226015,...
[3] 0.103757427817775,0.127648996987936,0.15250177579018,...
[4] 0.108024633674293,0.12118157243357,0.134813515447005,...
[5] 0.109476254160459,0.12659566766065,0.144359995969534,...
[6] 0.104671704031439,0.120304286111549,0.136527143979231,...

                                spline_der width_spline
                                <list>      <integer>
[1] 0.00823718792931595,0.00849663811461222,0.00875919479887256,...      593
[2] 0.0187191316937831,0.0194522839372348,0.020199731453538,...      335
[3] 0.0234173093333372,0.0243690129025921,0.025339543308228,...      508
[4] 0.0129217547880229,0.0133932818084443,0.0138717632963375,...      523
[5] 0.0168009286185096,0.0174398846432053,0.0180907582358953,...      315
[6] 0.0153411181741832,0.0159258829799662,0.0165216697493268,...      332
start_spline end_spline
<numeric>    <numeric>
[1] 3337483     3338075
[2] 4369075     4369409
[3] 4375417     4375924
[4] 4715698     4716220
[5] 4716303     4716617
[6] 4921234     4921565

[1] 0.00263609710663517,0.0030440063908373,...
[2] 0.00122834484766426,0.00147221033932346,0.00172653795382648,0.00199155049789865,0.002267...
[3] 0.000864564256647586,0.00120333309941394,0.00156495680464632,0.00195011992527804,0.002359...
[4] 0.000788310894250635,0.000956445833827378,0.00113504914292024,0.0013243784475671,0.001524...
[5] 0.00235269135248368,0.00273549871311195,0.00313329749625751,0.00354638183650247,0.00397503469483484,0.0043653050...
[6] 0.00139356234463668,0.0016220858943517,0.00186006713373021,0.00210769852286914,0.00236517387...

[1] 0.000238708662407768,0.000249059436995862,0.000259632908095442,0.000269751111111111,0.000279859162901525,0.0002899663157894737,0.000299999999999999,0.000309999999999999,0.000319999999999999,0.000329999999999999,0.000339999999999999,0.000349999999999999,0.000359999999999999,0.000369999999999999,0.000379999999999999,0.000389999999999999,0.000399999999999999,0.000409999999999999,0.000419999999999999,0.000429999999999999,0.000439999999999999,0.000449999999999999,0.000459999999999999,0.000469999999999999,0.000479999999999999,0.000489999999999999,0.000499999999999999,0.000509999999999999,0.000519999999999999,0.000529999999999999,0.000539999999999999,0.000549999999999999,0.000559999999999999,0.000569999999999999,0.000579999999999999,0.000589999999999999,0.000599999999999999,0.000609999999999999,0.000619999999999999,0.000629999999999999,0.000639999999999999,0.000649999999999999,0.000659999999999999,0.000669999999999999,0.000679999999999999,0.000689999999999999,0.000699999999999999,0.000709999999999999,0.000719999999999999,0.000729999999999999,0.000739999999999999,0.000749999999999999,0.000759999999999999,0.000769999999999999,0.000779999999999999,0.000789999999999999,0.000799999999999999,0.000809999999999999,0.000819999999999999,0.000829999999999999,0.000839999999999999,0.000849999999999999,0.000859999999999999,0.000869999999999999,0.000879999999999999,0.000889999999999999,0.000899999999999999,0.000909999999999999,0.000919999999999999,0.000929999999999999,0.000939999999999999,0.000949999999999999,0.000959999999999999,0.000969999999999999,0.000979999999999999,0.000989999999999999,0.000999999999999999,0.001009999999999999,0.001019999999999999,0.001029999999999999,0.001039999999999999,0.001049999999999999,0.001059999999999999,0.001069999999999999,0.001079999999999999,0.001089999999999999,0.001099999999999999,0.001109999999999999,0.001119999999999999,0.001129999999999999,0.001139999999999999,0.001149999999999999,0.001159999999999999,0.001169999999999999,0.001179999999999999,0.001189999999999999,0.001199999999999999,0.001209999999999999,0.001219999999999999,0.001229999999999999,0.001239999999999999,0.001249999999999999,0.001259999999999999,0.001269999999999999,0.001279999999999999,0.001289999999999999,0.001299999999999999,0.001309999999999999,0.001319999999999999,0.001329999999999999,0.001339999999999999,0.001349999999999999,0.001359999999999999,0.001369999999999999,0.001379999999999999,0.001389999999999999,0.001399999999999999,0.001409999999999999,0.001419999999999999,0.001429999999999999,0.001439999999999999,0.001449999999999999,0.001459999999999999,0.001469999999999999,0.001479999999999999,0.001489999999999999,0.001499999999999999,0.001509999999999999,0.001519999999999999,0.001529999999999999,0.001539999999999999,0.001549999999999999,0.001559999999999999,0.001569999999999999,0.001579999999999999,0.001589999999999999,0.001599999999999999,0.001609999999999999,0.001619999999999999,0.001629999999999999,0.001639999999999999,0.001649999999999999,0.001659999999999999,0.001669999999999999,0.001679999999999999,0.001689999999999999,0.001699999999999999,0.001709999999999999,0.001719999999999999,0.001729999999999999,0.001739999999999999,0.001749999999999999,0.001759999999999999,0.001769999999999999,0.001779999999999999,0.001789999999999999,0.001799999999999999,0.001809999999999999,0.001819999999999999,0.001829999999999999,0.001839999999999999,0.001849999999999999,0.001859999999999999,0.001869999999999999,0.001879999999999999,0.001889999999999999,0.001899999999999999,0.001909999999999999,0.001919999999999999,0.001929999999999999,0.001939999999999999,0.001949999999999999,0.001959999999999999,0.001969999999999999,0.001979999999999999,0.001989999999999999,0.001999999999999999,0.002009999999999999,0.002019999999999999,0.002029999999999999,0.002039999999999999,0.002049999999999999,0.002059999999999999,0.002069999999999999,0.002079999999999999,0.002089999999999999,0.002099999999999999,0.002109999999999999,0.002119999999999999,0.002129999999999999,0.002139999999999999,0.002149999999999999,0.002159999999999999,0.002169999999999999,0.002179999999999999,0.002189999999999999,0.002199999999999999,0.002209999999999999,0.002219999999999999,0.002229999999999999,0.002239999999999999,0.002249999999999999,0.002259999999999999,0.002269999999999999,0.002279999999999999,0.002289999999999999,0.002299999999999999,0.002309999999999999,0.002319999999999999,0.002329999999999999,0.002339999999999999,0.002349999999999999,0.002359999999999999,0.002369999999999999,0.002379999999999999,0.002389999999999999,0.002399999999999999,0.002409999999999999,0.002419999999999999,0.002429999999999999,0.002439999999999999,0.002449999999999999,0.002459999999999999,0.002469999999999999,0.002479999999999999,0.002489999999999999,0.002499999999999999,0.002509999999999999,0.002519999999999999,0.002529999999999999,0.002539999999999999,0.002549999999999999,0.002559999999999999,0.002569999999999999,0.002579999999999999,0.002589999999999999,0.002599999999999999,0.002609999999999999,0.002619999999999999,0.002629999999999999,0.002639999999999999,0.002649999999999999,0.002659999999999999,0.002669999999999999,0.002679999999999999,0.002689999999999999,0.002699999999999999,0.002709999999999999,0.002719999999999999,0.002729999999999999,0.002739999999999999,0.002749999999999999,0.002759999999999999,0.002769999999999999,0.002779999999999999,0.002789999999999999,0.002799999999999999,0.002809999999999999,0.002819999999999999,0.002829999999999999,0.002839999999999999,0.002849999999999999,0.002859999999999999,0.002869999999999999,0.002879999999999999,0.002889999999999999,0.002899999999999999,0.002909999999999999,0.002919999999999999,0.002929999999999999,0.002939999999999999,0.002949999999999999,0.002959999999999999,0.002969999999999999,0.002979999999999999,0.002989999999999999,0.002999999999999999,0.003009999999999999,0.003019999999999999,0.003029999999999999,0.003039999999999999,0.003049999999999999,0.003059999999999999,0.003069999999999999,0.003079999999999999,0.003089999999999999,0.003099999999999999,0.003109999999999999,0.003119999999999999,0.003129999999999999,0.003139999999999999,0.003149999999999999,0.003159999999999999,0.003169999999999999,0.003179999999999999,0.003189999999999999,0.003199999999999999,0.003209999999999999,0.003219999999999999,0.003229999999999999,0.003239999999999999,0.003249999999999999,0.003259999999999999,0.003269999999999999,0.003279999999999999,0.003289999999999999,0.003299999999999999,0.003309999999999999,0.003319999999999999,0.003329999999999999,0.003339999999999999,0.003349999999999999,0.003359999999999999,0.003369999999999999,0.003379999999999999,0.003389999999999999,0.003399999999999999,0.003409999999999999,0.003419999999999999,0.003429999999999999,0.003439999999999999,0.003449999999999999,0.003459999999999999,0.003469999999999999,0.003479999999999999,0.003489999999999999,0.003499999999999999,0.003509999999999999,0.003519999999999999,0.003529999999999999,0.003539999999999999,0.003549999999999999,0.003559999999999999,0.003569999999999999,0.003579999999999999,0.003589999999999999,0.003599999999999999,0.003609999999999999,0.003619999999999999,0.003629999999999999,0.003639999999999999,0.003649999999999999,0.003659999999999999,0.003669999999999999,0.003679999999999999,0.003689999999999999,0.003699999999999999,0.003709999999999999,0.003719999999999999,0.003729999999999999,0.003739999999999999,0.003749999999999999,0.003759999999999999,0.003769999999999999,0.003779999999999999,0.003789999999999999,0.003799999999999999,0.003809999999999999,0.003819999999999999,0.003829999999999999,0.003839999999999999,0.003849999999999999,0.003859999999999999,0.003869999999999999,0.003879999999999999,0.003889999999999999,0.003899999999999999,0.003909999999999999,0.003919999999999999,0.003929999999999999,0.003939999999999999,0.003949999999999999,0.003959999999999999,0.003969999999999999,0.003979999999999999,0.003989999999999999,0.003999999999999999,0.004009999999999999,0.004019999999999999,0.004029999999999999,0.004039999999999999,0.004049999999999999,0.004059999999999999,0.004069999999999999,0.004079999999999999,0.004089999999999999,0.004099999999999999,0.004109999999999999,0.004119999999999999,0.004129999999999999,0.004139999999999999,0.004149999999999999,0.004159999999999999,0.004169999999999999,0.004179999999999999,0.004189999999999999,0.004199999999999999,0.004209999999999999,0.004219999999999999,0.004229999999999999,0.004239999999999999,0.004249999999999999,0.004259999999999999,0.004269999999999999,0.004279999999999999,0.004289999999999999,0.004299999999999999,0.004309999999999999,0.004319999999999999,0.004329999999999999,0.004339999999999999,0.004349999999999999,0.004359999999999999,0.004369999999999999,0.004379999999999999,0.004389999999999999,0.004399999999999999,0.004409999999999999,0.004419999999999999,0.004429999999999999,0.004439999999999999,0.004449999999999999,0.004459999999999999,0.004469999999999999,0.004479999999999999,0.004489999999999999,0.004499999999999999,0.004509999999999999,0.004519999999999999,0.004529999999999999,0.004539999999999999,0.004549999999999999,0.004559999999999999,0.004569999999999999,0.004579999999999999,0.004589999999999999,0.004599999999999999,0.004609999999999999,0.004619999999999999,0.004629999999999999,0.004639999999999999,0.004649999999999999,0.004659999999999999,0.004669999999999999,0.004679999999999999,0.004689999999999999,0.004699999999999999,0.004709999999999999,0.004719999999999999,0.004729999999999999,0.004739999999999999,0.004749999999999999,0.004759999999999999,0.004769999999999999,0.004779999999999999,0.004789999999999999,0.004799999999999999,0.004809999999999999,0.004819999999999999,0.004829999999999999,0.004839999999999999,0.004849999999999999,0.004859999999999999,0.004869999999999999,0.004879999999999999,0.004889999999999999,0.004899999999999999,0.004909999999999999,0.004919999999999999,0.004929999999999999,0.004939999999999999,0.004949999999999999,0.004959999999999999,0.004969999999999999,0.004979999999999999,0.004989999999999999,0.004999999999999999,0.005009999999999999,0.005019999999999999,0.005029999999999999,0.005039999999999999,0.005049999999999999,0.005059999999999999,0.005069999999999999,0.005079999999999999,0.005089999999999999,0.005099999999999999,0.005109999999999999,0.005119999999999999,0.005129999999999999,0.005139999999999999,0.005149999999999999,0.005159999999999999,0.005169999999999999,0.005179999999999999,0.005189999999999999,0.005199999999999999,0.005209999999999999,0.005219999999999999,0.005229999999999999,0.005239999999999999,0.005249999999999999,0.005259999999999999,0.005269999999999999,0.005279999999999999,0.005289999999999999,0.005299999999999999,0.005309999999999999,0.005319999999999999,0.005329999999999999,0.00533999
```

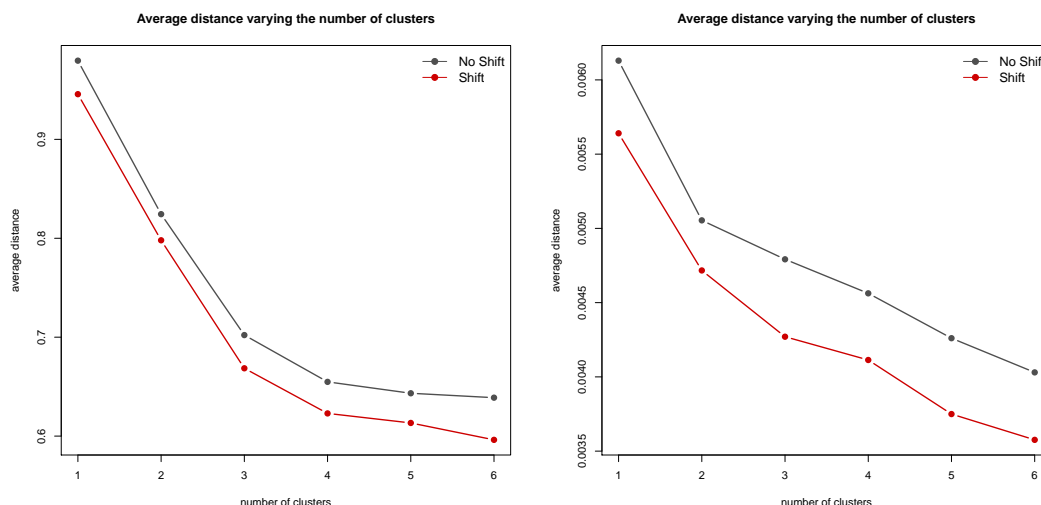


Figure 4: **Global distance within clusters.** Global distance of the peaks from the corresponding template, as a function of the number of clusters k . In the left panel the graph for the original spline approximation, while in the right panel results are relative to the scaled approximation.

```
[3] 0.00668285220715305,0.00627910073393782,0.00460131880636674,... 1,2,3,...
[4] 0.00461630755109992,0.00470065453190587,0.00449050254427536,... 1,2,2,...
[5] 0.00446076173955213,0.00382455339523424,0.00269385225273107,... 1,2,3,...
[6] 0.00130263542562362,0.00127949199373814,0.00132011451049066,... 1,2,2,...
      coef_shift                                dist_shift
      <list>                                <list>
[1] 1,36,41,... 0.0136019700774668,0.00529868829574966,0.00529868829574966,...
[2] 7,8,10,... 0.00288555205001011,0.00288555205001011,0.00196125806891698,...
[3] -17,-13,-4,... 0.00556533810032155,0.00561814086680531,0.00411761765869642,...
[4] -6,-2,-1,... 0.00469497672879888,0.00474230622233802,0.00488326800693634,...
[5] -13,-11,2,... 0.00330416023152333,0.00331770403459043,0.0057934253968632,...
[6] -3,-3,-2,... 0.00121592857872764,0.00123025283644485,0,...
```

```
-----
seqinfo: 20 sequences from an unspecified genome; no seqlengths
```

```
>
```

The particular case of k -mean alignment with $k = 1$ clusters can be used to highlight the effects of the alignment of the peaks: no grouping is performed, just the shifts are computed. Therefore, the decrease of the global distance is solely due to a change of the abscissae of the functions, as Figure ?? shows. Moreover, focusing for example on the first panel of Figure ??, we can deduce that, for this case

- the alignment can effectively decrease the distance, for example for $k = 6$, the gap between red and black line is significant;
- the alignment may change the optimal k : looking at the black line, one would have chosen $k = 4$, while the red line suggests $k = 3$ is the best choice. With the introduction of the shifts, data which are originally different becomes more similar and therefore one less cluster is needed; it has to be noted that the distance obtained with $k = 3$ and alignment is very similar to the one obtained with $k = 4$ and no alignment.

Therefore, for this case, one possible classification is the one associated to $k = 3$ with shift. On the contrary for the scaled peaks the value of k we can identify as crucial is $k = 2$ and shift is relevant since it reduces a lot the global distance. The results for this specific number of clusters can then be selected with the `choose_k` method:

```
> # select the results for k = 3 with alignment
```



```
> peaks.classified.short <- choose_k(peaks.cluster, k = 3,
+                                   shift = TRUE, cleaning = TRUE)
> head(peaks.classified.short)
```

GRanges object with 6 ranges and 1 metadata column:

| | seqnames | ranges | strand | cluster |
|-----|----------|--------------------|--------|-----------|
| | <Rle> | <IRanges> | <Rle> | <numeric> |
| [1] | chr18 | [3337524, 3338025] | * | 1 |
| [2] | chr18 | [4369126, 4369352] | * | 2 |
| [3] | chr18 | [4375448, 4375883] | * | 2 |
| [4] | chr18 | [4715744, 4716162] | * | 2 |
| [5] | chr18 | [4716374, 4716597] | * | 1 |
| [6] | chr18 | [4921270, 4921506] | * | 2 |

seqinfo: 20 sequences from an unspecified genome; no seqlengths

```
> peaks.classified.extended <- choose_k(peaks.cluster, k = 3,
+                                       shift = TRUE, cleaning = FALSE)
> # and for the scaled version for k =2 and alignment
>
> peaks.classified.scaled.short <- choose_k(peaks.cluster.scaled, k = 2,
+                                           shift = TRUE, cleaning = TRUE)
> head(peaks.classified.scaled.short)
```

GRanges object with 6 ranges and 4 metadata columns:

| | seqnames | ranges | strand | |
|-----|----------|--------------------|--------|--|
| | <Rle> | <IRanges> | <Rle> | |
| [1] | chr18 | [3337524, 3338025] | * | |
| [2] | chr18 | [4369126, 4369352] | * | |
| [3] | chr18 | [4375448, 4375883] | * | |
| [4] | chr18 | [4715744, 4716162] | * | |
| [5] | chr18 | [4716374, 4716597] | * | |
| [6] | chr18 | [4921270, 4921506] | * | |

| | |
|-----|--|
| [1] | 0.00263609710663517,0.0030440063908373 |
| [2] | 0.00122834484766426,0.00147221033932346,0.00172653795382648,0.00199155049789865,0.002267 |
| [3] | 0.000864564256647586,0.0012033309941394,0.00156495680464632,0.00195011992527804,0.002359 |
| [4] | 0.000788310894250635,0.000956445833827378,0.00113504914292024,0.0013243784475671,0.001524 |
| [5] | 0.00235269135248368,0.00273549871311195,0.00313329749625751,0.00354638183650247,0.00397503469483484, |
| [6] | 0.00139356234463668,0.0016220858943517,0.00186006713373021,0.00210769852286914,0.00236517387 |

| | |
|-----|--|
| [1] | |
| [2] | 0.000238708662407768,0.000249059436995862,0.000259632908095442,0.000269852286914,0.00027997503469483484, |
| [3] | 0.000327571761954349,0.000350081098788865,0.000373281486886386,0.000397503469483484,0.00042774124776864, |
| [4] | 0.000162985979570835,0.0001733265117332651,0.000183650247,0.00019469042, |
| [5] | 0.000375410625363164,0.00039025358389014,0.00040539347039775,0.00042082190469042,0.000436530506 |
| [6] | 0.000223858745449888,0.000233220374263462,0.000242774124776864,0.000252521012181378,0.000262462051668 |

| | summit_spline_rescaled | cluster |
|-----|------------------------|-----------|
| | <integer> | <numeric> |
| [1] | 227 | 1 |
| [2] | 168 | 2 |
| [3] | 104 | 2 |
| [4] | 180 | 2 |
| [5] | 117 | 2 |

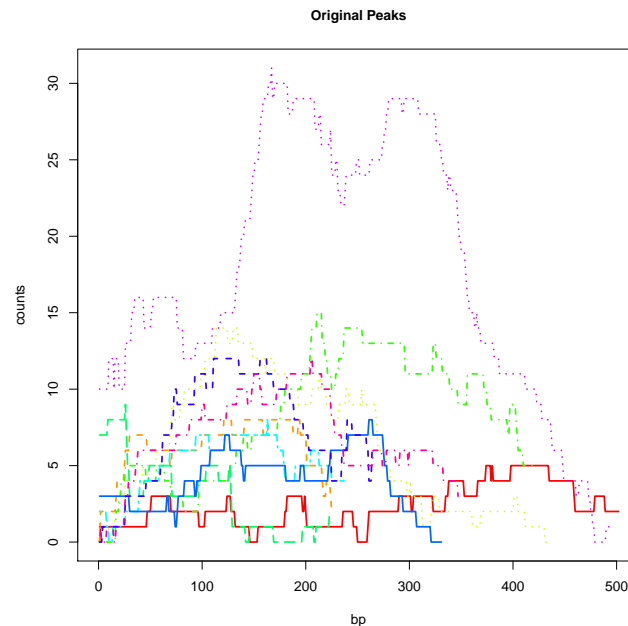


Figure 5: **10 peaks: counts.** Representation of the original peaks as raw counts (no smoothing).

```
[6]                155                2
```

```
-----
```

```
seqinfo: 20 sequences from an unspecified genome; no seqlengths
```

```
> peaks.classified.scaled.extended <- choose_k(peaks.cluster.scaled, k = 2,
+                                             shift = TRUE, cleaning = FALSE)
```

The `choose_k` method allows, respectively, to remove all the metadata columns computed by *FunChIP* and obtain a *GRanges* equivalent to the initial one, with an extra the metadata column `cluster` containing the classification labels (`cleaning = TRUE`), or a *GRanges* retaining all the details of the preprocessing and clustering (all the previously described metadata columns), with the extra column `cluster` (`cleaning = FALSE`).

5 Visualization of the peaks

The `plot_peak` method is a very flexible function for displaying ChIP-Seq peaks. In particular, it allows to plot the raw counts obtained by the `pileup_peak` method, as in Figure ???. It can also plot smoothed peaks, possibly centered around the summit, as in Figure ???, or scaled as in Figure ??? and centered.

```
> # plot of the first 10 peaks (raw data)
> plot_peak(peaks, index = 1:10, line.plot = 'count')
>
> # plot of the smoothed approximation of the first 10 peaks
> plot_peak(peaks.smooth, index = 1:10, line.plot = 'spline')
>
> # plot of the smoothed approximation of the first 10 peaks,
> # centering peaks around their summits
> plot_peak(peaks.summit, index = 1:10, line.plot = 'spline')
>
```

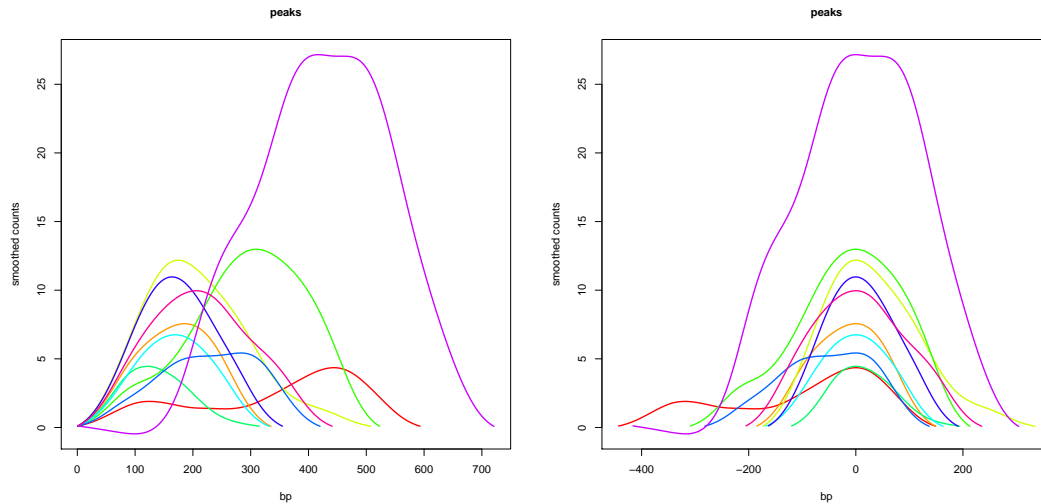


Figure 6: **10 spline-smoothed peaks.** In the left panel, smoothed peaks are shown, while in the right panel the same peaks are centered around their summits.

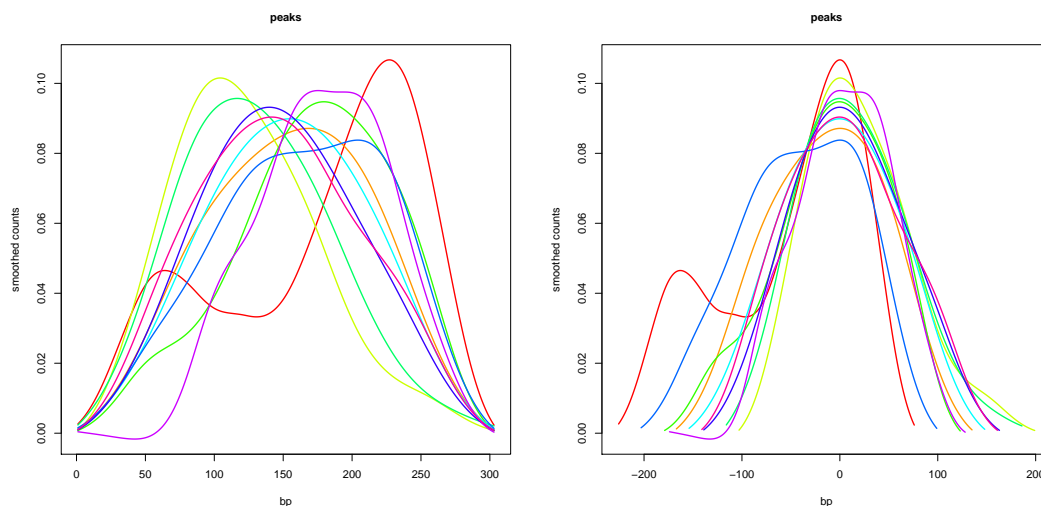


Figure 7: **10 spline-smoothed and scaled peaks.** In the left panel, scaled peaks are shown, while in the right panel the same peaks are centered around their summits.

```
> # plot of the smoothed approximation of the first 10 peaks;
> # the scaled functions are plotted.
> plot_peak(peaks.smooth.scaled, index = 1:10,
+           line.plot = 'spline', rescale = TRUE)
>

> # plot of the scaled approximation of the first 10 peaks,
> # centering peaks around their summits
> plot_peak(peaks.summit.scaled, index = 1:10,
+           line.plot = 'spline', rescale = TRUE)
>
```

From the comparison of Figure ?? and Figure ?? it is clear how the scaling affects the shape of splines. Now peaks are

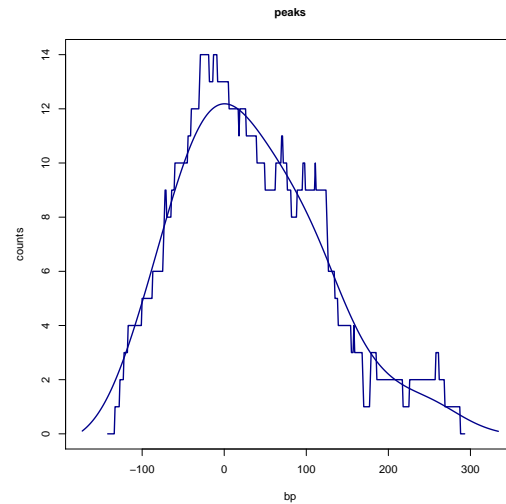


Figure 8: **Read coverage and spline approximation.** Plot of the original read coverage of a peak and its smoothing (spline approximation), centered around the summit.

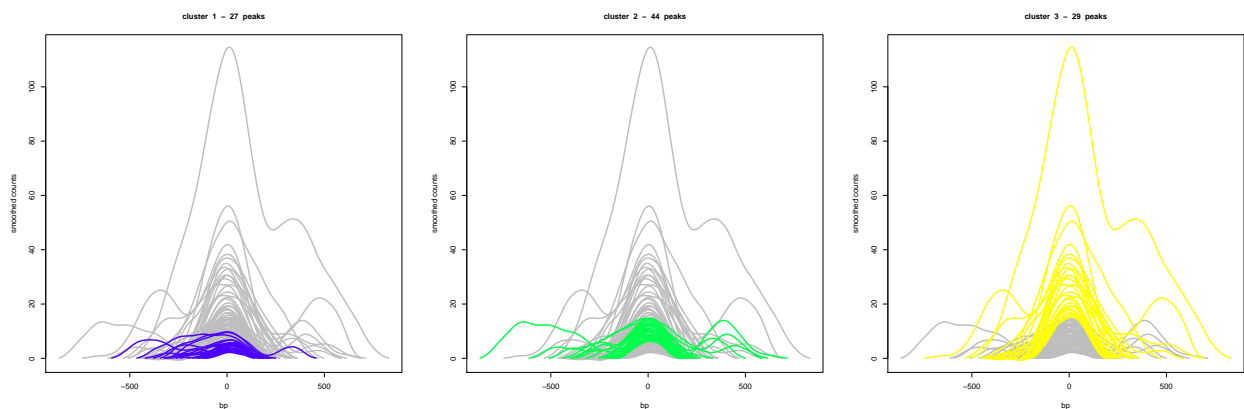


Figure 9: **Peaks divided in the three clusters** The same spline-smoothed peaks are plotted in grey, and for each panel the peaks in the corresponding cluster are colored to show their different shapes. Peaks are aligned with the shift coefficients obtained by the k-mean alignment algorithm.

no more related to the magnitude, but just to their shapes.

Moreover, plotting both raw counts and spline is also possible: Figure ?? shows a single peak in its raw and smoothed version. This representation is useful to check the accuracy of the smoothing and, if needed, manually set the λ parameter of the spline approximation.

```
> # plot of a peak comparing its raw structure and
> # its spline-smoothed version.
> plot_peak(peaks.summit, index = 3, line.plot = 'both', col = 'darkblue')
>
> # plot of the results of the kmean alignment.
> # Peaks are plotted in three different panels
> # according to the clustering results.
>
> plot_peak(peaks.cluster, index = 1:100, line.plot = 'spline',
```

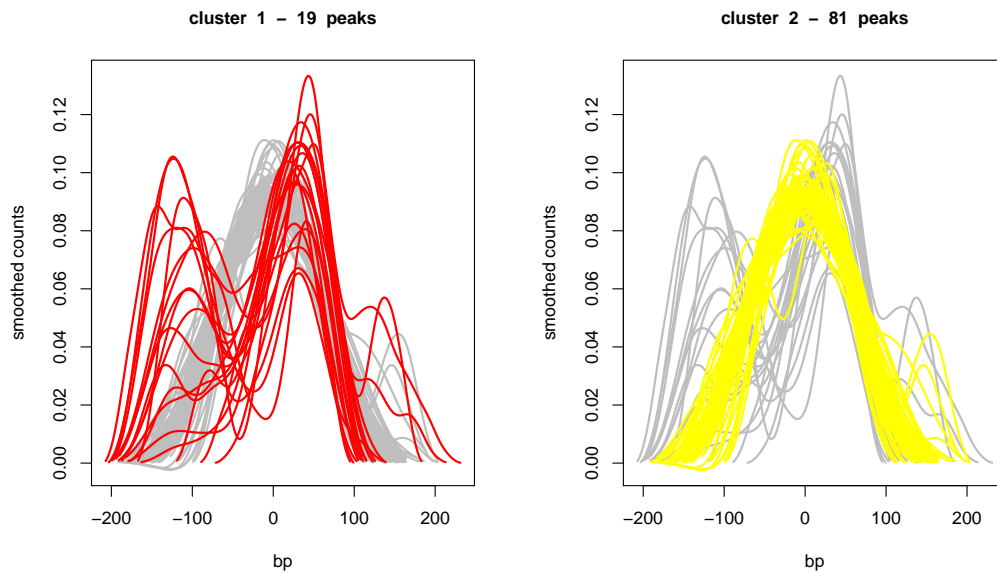


Figure 10: **Scaled peaks divided in the three clusters** The same spline-smoothed scaled peaks are plotted in grey, and for each panel the peaks in the corresponding cluster are colored to show their different shapes. Peaks are aligned with the shift coefficients obtained by the k-mean alignment algorithm.

```
+      shift = TRUE, k = 3, cluster.peak = TRUE,
+      col = topo.colors(3))
>
> # plot of the results of the kmean alignment.
> # Scaled peaks are plotted in three different panels
> # according to the clustering results.
>
> plot_peak(peaks.cluster.scaled, index = 1:100, line.plot = 'spline',
+          shift = TRUE, k = 2, cluster.peak = TRUE, rescale = TRUE,
+          col = heat.colors(2))
>
```

Finally, the `plot_peak` method allows to plot the results of the clustering via the k-mean alignment. In Figure ?? and Figure ??, smoothed and scaled peaks are divided into the three clusters and plotted with the optimal shift obtained with the alignment.