

Quick analysis of nucleosome positioning experiments using the **nucleR** package

Oscar Flores Guri

Institute for Research in Biomedicine & Barcelona Supercomputing Center
Joint Program on Computational Biology

October 1, 2012

1 Introduction

The **nucleR** package provides a high-level processing of genomic datasets focused in nucleosome positioning experiments, despite they should be also applicable to chromatin immunoprecipitation (ChIP) experiments in general.

The aim of this package is not providing an all-in-one data analysis pipeline but complement those existing specialized libraries for low-level data importation and pre-processing into R/Bioconductor framework.

nucleR works with data from the two main high-throughput technologies available nowadays for ChIP: Next Generation Sequencing/NGS (ChIP-seq) and Tiling Microarrays (ChIP-on-Chip).

This is a brief summary of the main functions:

- Data importation: `processReads`, `processTilingArray`
- Data transformation: `coverage.rpm`, `filterFFT`, `controlCorrection`
- Nucleosome calling: `peakDetection`, `peakScoring`
- Visualization: `plotPeaks`
- Data generation: `syntheticNucMap`

For more details about the functions and how to use them refer to the **nucleR** manual.

This software was published in Bioinformatics Journal. See the paper for additional information.[1]

2 Reading data

As mentioned previously, **nucleR** uses the pre-processed data of other lower level packages for data importation, supporting a few but common formats that should fulfill the requirements of most users..

`ExpressionSet` from package `Biobase` [2] is used for Tiling Array experiments as described in `Starr` [3] and other packages for the Tiling Array manipulation. This kind of experiments can be readed with the `processTilingArray` function.

`AlignedRead` from package `ShortRead` [4] is recommended for NGS, covering most of the state of the art sequencing technologies. Additionally, support for reads in `RangedData` format is also provided (a range per read with a `"strand"` column) .

2.1 Reading Tiling Arrays

Tiling Arrays are a cheap and fast way to have low-resolution nucleosome coverage maps. They have been widely used in literature[5, 6, 7], but complex statistical methods were needed for their processing [8].

This kind of microarrays cover a part of the genome with certain spacing between probes which causes a drop in the resolution and originates some problems. The nucleosome calling from Tiling Array data required hard work on bioinformatics side and use of heavy and artificial statistical machinery such as Hidden Markov Models [5, 6] or higher order Bayesian Networks [9].

`nucleR` presents a new method based on a simple but effective peak calling method which achieves a great performance at low computing cost that will be presented in subsequent sections.

In order to standardize the data coming both from Tiling Arrays and NGS, the array fluorescence intensities (usually the ratio of the hybridization of nucleosomal and control sample) are converted to 1bp resolution by inferring the missed values from the neighboring probes. This is done by the function `processTilingArray`:

```
processTilingArray(data, exprName, chrPattern, inferLen=50)
```

An example of a processed dataset is provided in this package. See the help page of `tilingArray_preproc` for details on how it has been created. This object is a numeric vector covering the 8000 first positions of chromosome 1 in yeast (*Saccharomices Cerevisiae* genome (SacCer1)).

```
> require(IRanges)
> library(nucleR)
> data(nucleosome_tiling)
> head(nucleosome_tiling, n=25)
```

```
[1] 1.273222 1.281978 1.290734 1.299490 1.308246 1.352696 1.397145 1.441595
[9] 1.486044 1.501795 1.517547 1.533298 1.549049 1.547577 1.546105 1.544633
[17] 1.543161 1.539886 1.536612 1.533337 1.530063 1.488922 1.447782 1.406642
[25] 1.365502
```

This values represent the normalized fluorescence intensity from hybridized sample of nucleosomal DNA versus naked DNA obtained from `Starr`. The values can be either direct observations (if a probe was starting at that position) or a inferred value from neighboring probes. This data can be passed directly to the filtering functions, as described later in the section 3.

2.2 Next Generation Sequencing

NGS has become one of the most popular technique to map nucleosome in the genome in the last years [10, 11, 12]. The drop of the costs of a genome wide sequencing together with the high resolution coverage maps obtained, made it the election of many scientists.

The package `ShortRead` allows reading of the data coming from many sources (Bowtie, MAQ, Illumina pipeline...) and has become one of the most popular packages in R/Bioconductor for NGS data manipulation.

A new R package, called `htSeqTools` [13], has been recently created to perform preprocessing and quality assesment on NGS experiments. `nucleR` supports most of the output generated by the functions on that package and recommends its use for quality control and correction of common biases that affect NGS.

`nucleR` handles `ShortRead` and `RangedData` data formats. The dataset `nucleosome_htseq` includes some NGS reads obtained from a nucleosome positioning experiment also from yeast genome, following a protocol similar to the one described in [6].

The paired-end reads coming from Illumina Genome Analyzer II sequencer were mapped using Bowtie and imported into R using `ShortRead`. Paired ends where merged and sorted according the start position. Those in the first 8000bp of chromosome 1 where saved for this example. Further details are in the reference [14]:

```
> data(nucleosome_htseq)
> class(nucleosome_htseq)
```

```
[1] "RangedData"
attr(,"package")
[1] "IRanges"
```

```
> nucleosome_htseq
```

RangedData with 18001 rows and 1 value column across 1 space

	space	ranges		strand
	<factor>	<IRanges>		<character>
1	chr1	[1, 284]		+
2	chr1	[5, 205]		+
3	chr1	[5, 205]		+
4	chr1	[5, 209]		+
5	chr1	[5, 283]		+
6	chr1	[6, 285]		+
7	chr1	[6, 285]		+
8	chr1	[8, 126]		+
9	chr1	[8, 127]		+
...
17993	chr1	[7994, 8148]		+
17994	chr1	[7994, 8150]		+
17995	chr1	[7994, 8150]		+

17996	chr1	[7994, 8151]		+
17997	chr1	[7994, 8151]		+
17998	chr1	[7994, 8151]		+
17999	chr1	[7994, 8151]		+
18000	chr1	[7994, 8152]		+
18001	chr1	[7994, 8152]		+

Now we will transform the reads to a normalized format. Moreover, as the data is paired-ended and we are only interested in mononucleosomes (which are typically 147bp), we will discard the reads with a length greater than 200bp, allowing margin for some underdigestion but discarding extra long reads. Note that the behaviour of `fragmentLen` is different for single-ended data, see the manual page of this function for detailed information.

As our final objective is identifying the nucleosome positions, and `nucleR` does it from the dyad, we will increase the sharpness of the dyads by removing some bases from the ends of each read. In the next example, we will create two new objects, one with the original paired-end reads and another one with the reads trimmed to the middle 40bp around the dyad (using the `trim` argument).

```
> #Process the paired end reads, but discard those with length > 200
> reads_pair = processReads(nucleosome_htseq, type="paired", fragmentLen=200)
> #Process the reads, but now trim each read to 40bp around the dyad
> reads_trim = processReads(nucleosome_htseq, type="paired", fragmentLen=200, trim=40)
```

The next step is obtain the coverage (the count of how many reads are in each position). The standard `IRanges` package function `coverage` will work well here, but it is a common practice to normalize the coverage values according to the total number of short reads obtained in the NGS experiment. The common used unit is *reads per milion* (r.p.m.) which is the coverage value divided by the total number of reads and multiplied per one milion. A quick and efficient way to do this with `nucleR` is the `coverage.rpm` function.¹

```
> #Calculate the coverage, directly in reads per million (r.p.m)
> cover_pair = coverage.rpm(reads_pair)
> cover_trim = coverage.rpm(reads_trim)
```

In Figure 1 we can observe the effect of `trim` attribute plotting both coverages. Note that the coverages are normalized in the range 0–1:

2.3 MNase bias correction

The Microccocal Nuclease is a widely used enzyme that has been proved to have a biase for certain dinucleotide steps [14]. In this package we offer a quick way to inspect the effect of such artifact by correcting the profiles of nucleosomal DNA reads with a mock sample of naked DNA digested with MNase.

¹Note that conversion in the example dataset gives huge values. This is because r.p.m. expects a large number of reads, and this dataset is only a fraction of a whole one. Also take into account that reads from single-ended (or trimmed reads) and reads from paired-ended could have different mean value of coverage

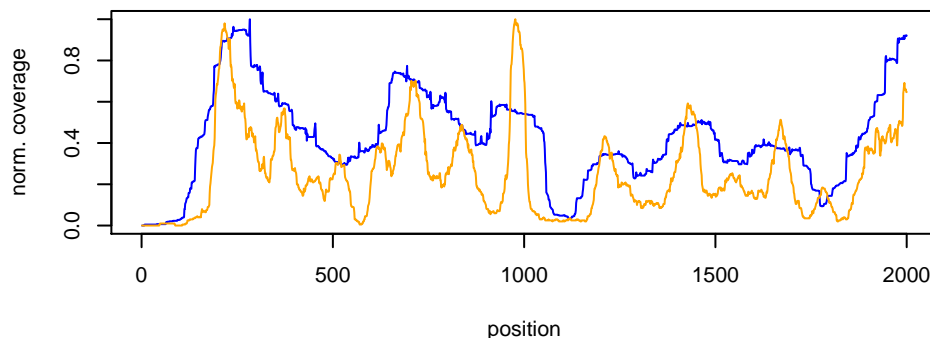


Figure 1: Variation in the sharpness of the peaks using `trim` attribute. In blue, the original coverage; in orange the trimmed version

The use of this function requires a paired-end control sample and a paired end or extended single-read nucleosomal DNA sample. A toy example generated using synthetic data can be found in Figure 2.

```
> #Toy example
> map = syntheticNucMap(as.ratio=TRUE, wp.num=50, fuz.num=25)
> exp = coverage(map$syn.reads)
> ctr = coverage(map$ctr.reads)
> corrected = controlCorrection(exp, ctr)
```

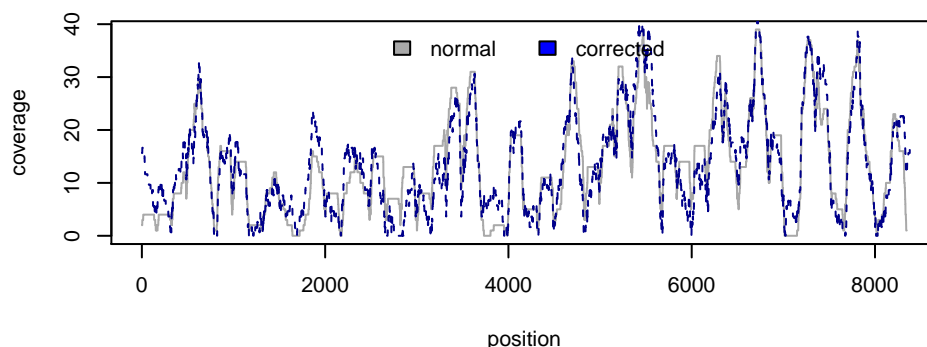


Figure 2: Toy example of MNase bias correction. Random nucleosomal and control reads have been generated using `syntheticNucMap` function and corrected using `controlCorrect`

3 Signal Smoothing and Nucleosome Calling

In the previous sections we converted the experimental data from NGS or Tiling Arrays to a continuous, 1bp resolution signal. In this section we will remove the noise present in the data and score the peaks identified, giving place to the nucleosome calls.

Previously, in the literature, Hidden Markov Models, Support Vector Machines or other complex intelligent agents were used for this task [5, 6, 9, 15, 12]. This was needed for dealing with the noise and uncertain characterization of the fuzzy positioning of the nucleosomes.

Despite this approach is a valid way to face the problem, the use of such artificial constructs is difficult to implement and sometimes requires a subjective modeling of the solution, constraining or at least conditioning the results observed.

The method presented here proposes to *keep it simple*, allowing the researcher to study the results he or she is interested *a posteriori*.

nucleR aim is to evaluate where the nucleosomes are located and how accurate that position is. We can find a nucleosome read in virtually any place in the genome, but some positions will show a high concentration and will allow us to mark this nucleosome as **well-positioned** whereas other will be less phased giving place to **fuzzy** or **de-localized** nucleosomes [16].

We think it's better to provide a detailed but convenient identification of the relevant nucleosome regions and score them according to its grade of fuzziness. From our point of view, every researcher should make the final decision regarding filtering, merging or classifying the nucleosomes according its necessities, and nucleR is only a tool to help in this "dirty" part of the research.

3.1 Noise removal

NGS and specially Tiling Array data show a very noisy profile which complicates the process of the nucleosome detection from peaks in the signal. A common approach used in the literature is smooth the signal with a sliding window average and then use a Hidden Markov Model to calculate the probabilities of having one or another state.

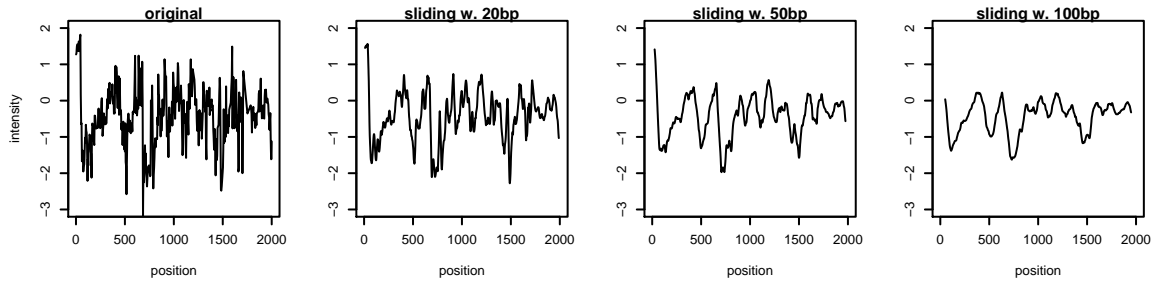


Figure 3: Original intensities from tiling array experiment. Smoothing using a sliding window of variable length (0, 20, 50 and 100 bp) is presented.

As can be seen in Figure 3, data needs some smoothing to be interpretable, but a simple sliding window average is not sufficient. Short windows allow too much noise but larger ones change the position and the shape of the peaks.

nucleR proposes a method of filtering based on the Fourier Analysis of the signal and the selection of its principal components.

Any signal can be described as a function of individual periodic waves with different frequencies and the combination of them creates more complex signals. The noise in a signal can be described as a small, non periodic fluctuations, and can be easily identified and removed [17].

nucleR uses this theory to transform the input data into the Fourier space using the Fast Fourier Transform (FFT). A FFT has a real and a imaginary component. The representation of the real component it's called the power spectrum of the signals and shows which are the frequencies that have more weight (power) in the signal. The low frequency components (so, very periodic) usually have a huge influence in the composite signal, but its relevance drops as

the frequency increases.

We can look at the power spectrum of the example dataset with the following command:

```
> fft_ta = filterFFT(nucleosome_tiling, pcKeepComp=0.01, showPowerSpec=TRUE)
```

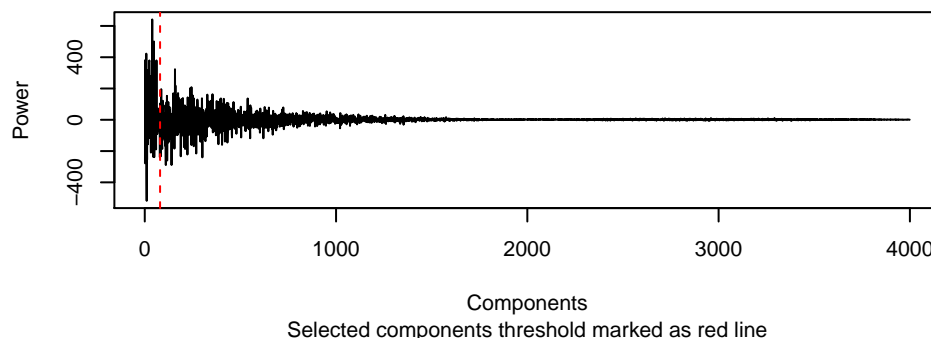


Figure 4: Power spectrum of the example Tiling Array data, percentile 1 marked with a dashed line

In the Figure 4 only the half of the components are plotted, as the spectrum is repeated symmetrically respect to its middle point. The first component (not shown in the plot), has period 1, and, in practice, is a count of the length of the signal, so it has a large value.

High frequency signals are usually echoes (repeating waves) of lower frequencies, i.e. a peak at 10 will be the sum of the pure frequency 10 plus the echo of the frequency 5 in its 2nd repetition. Echoes can be ignored without losing relevant information.

The approach **nucleR** follows is supposing that with just a small percentage of the components of the signal, the input signal can be recreated with a high precision, but without a significant amount of noise. We check empirically that with 1% or 2% of the components (this means account 1 or 2 components for each 100 positions of the genomic data) it's enough to recreate the signal with a very high correlation (>0.99). Tiling Array could require more smoothing (about 1% should be fine) and NGS has less noise and more components can be selected for fitting better the data (about 2%), See Figure 4 for the selected components in the example.

In order to easy the choice of the **pcKeepComp** parameter, **nucleR** includes a function for automatic detection of a fitted value that provides a correlation between the original and the filtered profiles close to the one specified. See the manual page of **pcKeepCompDetect** for detailed information.

In short, the cleaning process consists on converting the coverage/intensity values to the Fourier space, and knock-out (set to 0) the components greater than the given percentile in order to remove the noise from the profile. Then the inverse Fast Fourier Transform is applied to recreate the filtered signal. In Figure 5 the filtered signal is overlapped to the raw signal.

The cleaning of the input has almost no effect on the position and shape of the peaks, maintaining a high correlation with the original signal but allowing achieve a great performance with a simple peak detection algorithm:

```
> tiling_raw = nucleosome_tiling
> tiling_fft = filterFFT(tiling_raw, pcKeepComp=0.01)
```

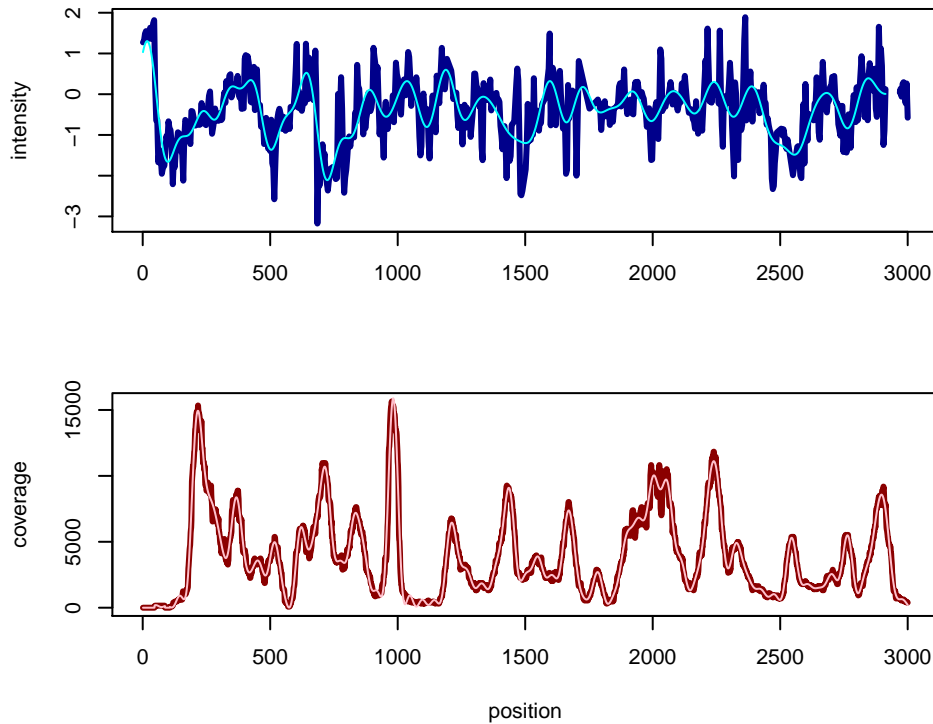


Figure 5: Filtering in Tiling Array (up, blue) (1% comp.) and NGS (down, red) (2%comp.)

```
> htseq_raw = as.vector(cover_trim[[1]])
> htseq_fft = filterFFT(htseq_raw, pcKeepComp=0.02)
> cor(tiling_raw, tiling_fft, use="complete.obs")
```

```
[1] 0.7153782
```

```
> cor(htseq_raw, htseq_fft, use="complete.obs")
```

```
[1] 0.9937643
```

3.2 Peak detection and Nucleosome Calling

After noise removal, the calling for nucleosomes is easy to perform. In nucleosome positioning, in contrast with other similar experiments like ChIP, the problem for the peaks detection algorithms is deal with the presence of an irregular signal which causes lots of local maxima (i.e., peaks due to noise inside a real peak). Here, we avoid this problem applying the FFT filter, allowing the detection of peaks in a simple but efficient way just looking for changes in the trend of the profile. This is implemented in the `peakDetection` function and results can be represented with the function `plotPeaks`:

```
> peaks = peakDetection(htseq_fft, threshold="25%", score=FALSE)
> peaks
```



```
[1] 218 368 452 518 623 715 776 836 983 1213 1331 1437 1549 1603 1672
[16] 1785 1945 2005 2053 2241 2330 2546 2605 2702 2764 2899 3124 3285 3354 3518
[31] 3710 3869 4009 4137 4233 4305 4383 4524 4596 4832 4912 4981 5171 5241 5291
[46] 5366 5458 5529 5598 5688 5752 5906 5978 6065 6123 6238 6312 6410 6472 6669
[61] 6834 6892 7002 7048 7100 7152 7314 7405 7457 7545 7615 7684 7775 7932 8042
```

```
> plotPeaks(peaks, htseq_fft, threshold="25%")
```

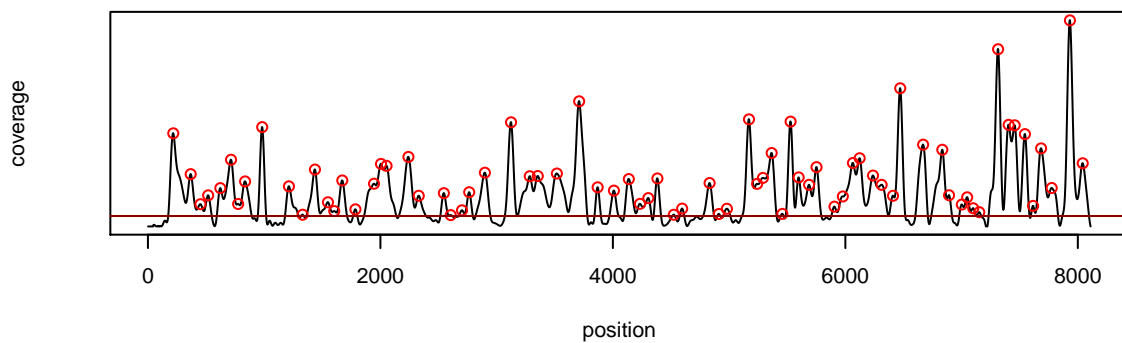


Figure 6: Output of `plotPeaks` function. Peaks are spotted in red and detection threshold marked with an horizontal line.

All the peaks above a threshold value are identified. Threshold can be set to 0 for detecting all the peaks, but this is not recommended as usually small fluctuations can appear in bottom part of the profile. This package also provides an automatic scoring of the peaks, which accounts for the two main features we are interested in: the height and the sharpness of the peak.

The *height* of a peak is a direct measure of the reads coverage in the peak position, but represented as a probability inside a Normal distribution.

The *sharpness* is a measure of how fuzzy is a nucleosome. If a peak is very narrow and the surrounding regions are depleted, this is an indicator of a good positioned nucleosome, while wide peaks or peaks very close to each other are probably fuzzy nucleosomes (despite the coverage can be very high in this region).

Scores can be calculated with the `peakScoring` function or directly with the argument `score=TRUE` in `peakDetection`.

```
> peaks = peakDetection(htseq_fft, threshold="25%", score=TRUE)
> head(peaks)
```

	peak	score
1	218	0.9749612
2	368	0.6824909
3	452	0.2675856
4	518	0.3829457
5	623	0.4858367
6	715	0.8408881

The scores in Figure 7 only account for the punctual height of the peak. As said previously, this measure can be improved by accounting the fuzzyness of a nucleosome call (the sharpness

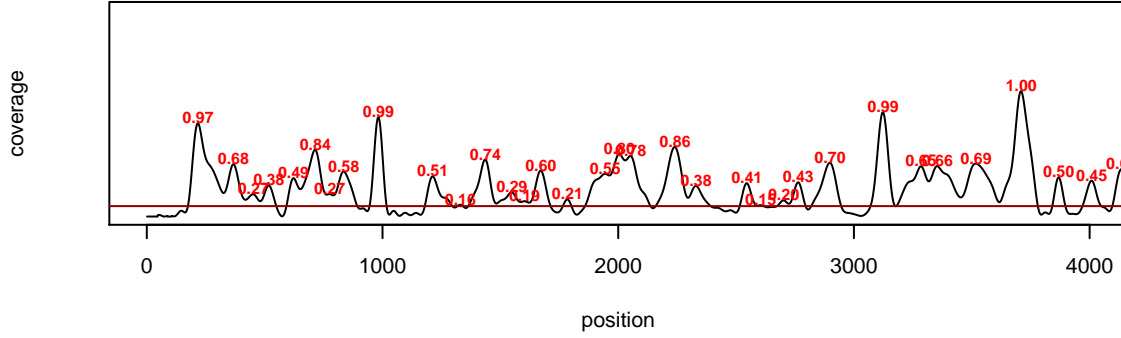


Figure 7: `plotPeaks` function with `score=TRUE`.

of the peak). This requires a way to account for longer range peaks, which can be obtained with the `width` argument. In this way one can convert the identified nucleosome dyads to whole nucleosome length ranges and account for its degree of fuzzyness:

```
> peaks = peakDetection(htseq_fft, threshold="25%", score=TRUE, width=140)
> head(peaks)
```

RangedData with 6 rows and 3 value columns across 1 space

	space	ranges	score	score_w	score_h
	<factor>	<IRanges>	<numeric>	<numeric>	<numeric>
1	1	[148, 287]	0.8213296	0.6676980	0.9749612
2	1	[298, 437]	0.6385168	0.5945427	0.6824909
3	1	[382, 521]	0.3157506	0.3639156	0.2675856
4	1	[448, 587]	0.5182733	0.6536008	0.3829457
5	1	[553, 692]	0.5182437	0.5506507	0.4858367
6	1	[645, 784]	0.7229812	0.6050742	0.8408881

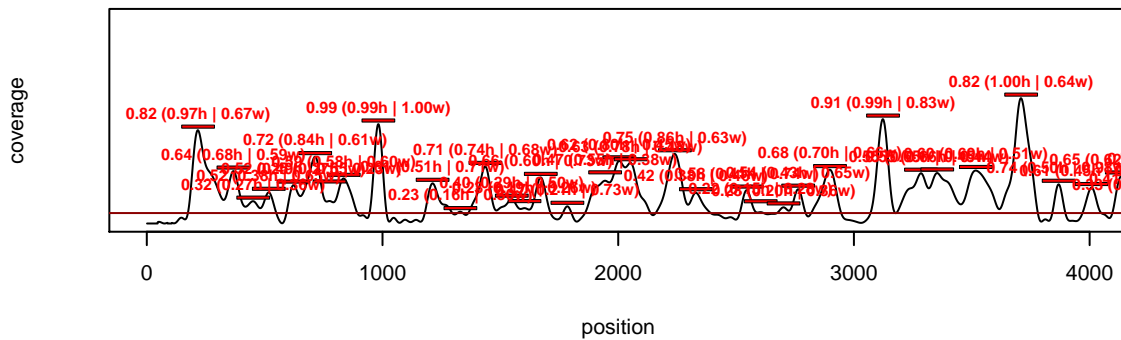


Figure 8: `plotPeaks` output with `score=TRUE` and `width=140`.

Note than in Figure 8 overlapped peaks in a width and tall region are penalized, meanwhile the peaks with surrounding depleted regions have a higher relative score. This is the approach recommended for working with nucleosome calls.

Nucleosome calls filtering, merging or classification can be performed with standard `IRanges` [18] functions, such as `reduce`, `findOverlaps` or `disjoint`.

The next example shows a simple way to merge those nucleosomes which are overlap accounting them as a fuzzy regions:

```
> nuc_calls = ranges(peaks[peaks$score > 0.1,])[[1]]
> red_calls = reduce(nuc_calls)
> red_class = RangedData(red_calls, isFuzzy=width(red_calls) > 140)
> head(red_class)
```

RangedData with 6 rows and 1 value column across 1 space

	space	ranges	isFuzzy
	<factor>	<IRanges>	<logical>
1	1	[148, 287]	FALSE
2	1	[298, 905]	TRUE
3	1	[913, 1052]	FALSE
4	1	[1143, 1854]	TRUE
5	1	[1875, 2122]	TRUE
6	1	[2171, 2399]	TRUE

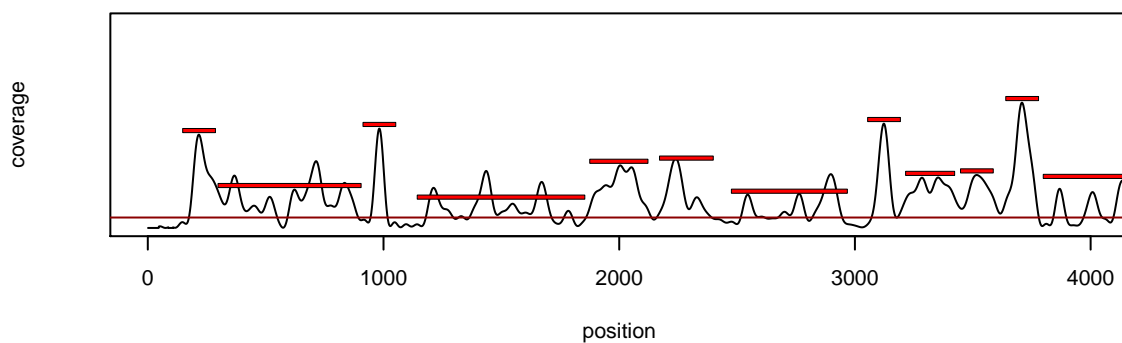


Figure 9: Simple example of ranges manipulation to plot fuzzy nucleosomes

4 Exporting data

`export.wig` and `export.bed` allow exportation of coverage/intensity values and nucleosome calls in a standard format which works on most of the genome browsers available today (like UCSC Genome Browser or Integrated Genome Browser).

`export.wig` creates WIG files which are suitable for coverage/intensities, meanwhile `export.bed` creates BED files which contain ranges and scores information, suitable for calls.

5 Generating synthetic maps

nucleR includes a synthetic nucleosome map generator, which can be helpful in benchmarking or comparing data against a random map. `syntheticNucMap` function does that, allowing a full customization of the generated maps.

When generating a map, the user can choose the number of the well-positioned and fuzzy nucleosome, as their variance or maximum number of reads. It also provides an option to calculate the ratio between the generated nucleosome map and a mock control of random reads (like a naked DNA randomly fragmented sample) to simulate hybridization data of Tiling Arrays.

The perfect information about the nucleosome dyads is returned by this function, together with the coverage or ratio profiles.

See the man page of this function for detailed information about the different parameters and options.

```
> syntheticNucMap(wp.num=100, wp.del=10, wp.var=30, fuz.num=20, fuz.var=50,
+ max.cover=20, nuc.len=147, lin.len=20, rnd.seed=1, as.ratio=TRUE, show.plot=TRUE)
```

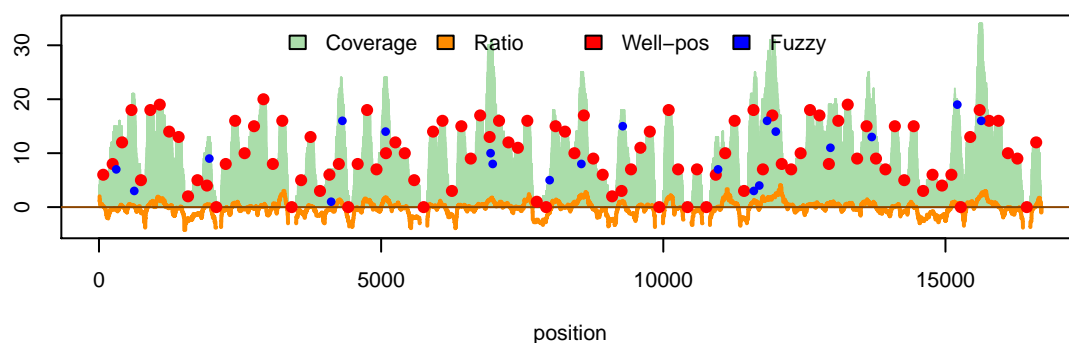


Figure 10: Example synthetic coverage map of 90 well-positioned (100-10) and 20 fuzzy nucleosomes.

References

- [1] O. Flores and M. Orozco. nucleR: a package for non-parametric nucleosome positioning. *Bioinformatics*, early publication online:–, 2011.
- [2] Robert C Gentleman, Vincent J. Carey, Douglas M. Bates, and others. Bioconductor: Open software development for computational biology and bioinformatics. *Genome Biology*, 5:R80, 2004.
- [3] Benedikt Zacher, Johannes Soeding, Pei Fen Kuan, Matthias Siebert, and Achim Tresch. *Starr: Simple tiling array analysis of Affymetrix ChIP-chip data*, 2009.
- [4] Martin Morgan, Simon Anders, Michael Lawrence, Patrick Aboyoun, Hervé Pagès, and Robert Gentleman. ShortRead: a Bioconductor package for input, quality assessment and exploration of high-throughput sequence data. *Bioinformatics*, 25:2607–2608, 2009.
- [5] Guo-Cheng Yuan, Yuen-Jong Liu, Michael F Dion, Michael D Slack, Lani F Wu, Steven J Altschuler, and Oliver J Rando. Genome-scale identification of nucleosome positions in *S. cerevisiae*. *Science (New York, N.Y.)*, 309:626–30, 2005.
- [6] William Lee, Desiree Tillo, Nicolas Bray, Randall H Morse, Ronald W Davis, Timothy R Hughes, and Corey Nislow. A high-resolution atlas of nucleosome occupancy in yeast. *Nature genetics*, 39:1235–44, 2007.

- [7] Travis N Mavrich, Cizhong Jiang, Ilya P Ioshikhes, Xiaoyong Li, Bryan J Venters, Sara J Zanton, Lynn P Tomsho, Ji Qi, Robert L Glaser, Stephan C Schuster, David S Gilmour, Istvan Albert, and B Franklin Pugh. Nucleosome organization in the *Drosophila* genome. *Nature*, 453:358–62, 2008.
- [8] X Shirley Liu. Getting started in tiling microarray analysis. *PLoS computational biology*, 3:1842–4, 2007.
- [9] Pei Fen Kuan, Dana Huebert, Audrey Gasch, and Sunduz Keles. A non-homogeneous hidden-state model on first order differences for automatic detection of nucleosome positions. *Statistical applications in genetics and molecular biology*, 8:Article29, 2009.
- [10] Noam Kaplan, Irene K Moore, Yvonne Fondufe-Mittendorf, Andrea J Gossett, Desiree Tillo, Yair Field, Emily M LeProust, Timothy R Hughes, Jason D Lieb, Jonathan Widom, and Eran Segal. The DNA-encoded nucleosome organization of a eukaryotic genome. *Nature*, 458:362–6, 2009.
- [11] Dustin E Schones, Kairong Cui, Suresh Cuddapah, Tae-Young Roh, Artem Barski, Zhibin Wang, Gang Wei, and Keji Zhao. Dynamic regulation of nucleosome positioning in the human genome. *Cell*, 132:887–98, 2008.
- [12] Liqun Xi, Yvonne Fondufe-Mittendorf, Lei Xia, Jared Flatow, Jonathan Widom, and Ji-Ping Wang. Predicting nucleosome positioning using a duration Hidden Markov Model. *BMC Bioinformatics*, 11:346, 2010.
- [13] David Rossell, Evarist Planet, Camille Stephan-Otto, and Oscar Flores. *htSeqTools: High-Throughput Sequencing Data Analysis*, 2011.
- [14] Özgen Deniz, Oscar Flores, Federica Battistini, Alberto Pérez, Montserrat Soler-López, and Modesto Orozco. Physical Properties of Naked DNA Signal Gene Regulatory Regions. *Work under submission*, 2011.
- [15] X. Chen, M. M. Hoffman, J. a. Bilmes, J. R. Hesselberth, and W. S. Noble. A dynamic Bayesian network for identifying protein-binding footprints from single molecule-based sequencing data. *Bioinformatics*, 26:i334–i342, 2010.
- [16] Cizhong Jiang and B Franklin Pugh. Nucleosome positioning and gene regulation: advances through genomics. *Nature reviews. Genetics*, 10(3):161–72, 2009.
- [17] Steven W. Smith. *The Scientist and Engineer’s Guide to Digital Signal Processing (Second ed.)*. California Technical Publishing, 1999. Available online: <http://www.dspguide.com/pdfbook.htm>.
- [18] H. Pages, P. Aboyoun, and M. Lawrence. *IRanges: Infrastructure for manipulating intervals on sequences*.