

# Analysis of ChIP-seq Data with ‘**mosaics**’ Package

Dongjun Chung<sup>1</sup>, Pei Fen Kuan<sup>2</sup> and Sündüz Keleş<sup>1,3</sup>

<sup>1</sup>Department of Statistics, University of Wisconsin  
Madison, WI 53706.

<sup>2</sup>Department of Biostatistics, University of North Carolina at Chapel Hill  
Chapel Hill, NC 27599.

<sup>3</sup>Department of Biostatistics and Medical Informatics, University of Wisconsin  
Madison, WI 53706.

June 21, 2012

## Contents

|   |           |
|---|-----------|
| <b>1 Overview</b>   | <b>1</b>  |
| <b>2 Getting started</b>  | <b>2</b>  |
| <b>3 Two-Sample Analysis using ‘mosaicsRunAll’</b>                    | <b>2</b>  |
| <b>4 Workflow: Two-Sample Analysis</b>                                | <b>3</b>  |
| 4.1 Constructing Bin-Level Files from the Aligned Read File . . . . . | 3         |
| 4.2 Reading Bin-Level Data into the R Environment . . . . .           | 5         |
| 4.3 Fitting the MOSAiCS Model . . . . .                               | 6         |
| 4.4 Identifying Peaks Based on the Fitted Model . . . . .             | 8         |
| <b>5 Two-Sample Analysis with Mappability and GC Content</b>          | <b>11</b> |
| <b>6 One-Sample Analysis</b>  | <b>13</b> |
| <b>7 Conclusion and Ongoing Work</b>                                  | <b>14</b> |

## 1 Overview

This vignette provides an introduction to the analysis of ChIP-seq data with the ‘**mosaics**’ package. R package **mosaics** implements MOSAiCS, a statistical framework for the analysis of ChIP-seq data, proposed in [1]. MOSAiCS stands for “**MO**del-based one and two **S**ample **A**nalysis and **I**nfERENCE for **ChIP-Seq** Data”. It implements a flexible parametric mixture modeling approach for detecting peaks, i.e., enriched regions, in one-sample (ChIP sample) or two-sample (ChIP and control samples) ChIP-seq data. It can account for mappability and GC content biases that arise in ChIP-seq data.

The package can be loaded with the command:

```
R> library("mosaics")
```

## 2 Getting started

'mosaics' package provides flexible framework for the ChIP-seq analysis.

If you have the data for matched control sample, two-sample analysis is recommended. If the ChIP-seq data is deeply sequenced, the two-sample analysis without mappability and GC content (Section 4) is usually appropriate. For the ChIP-seq data with low sequencing depth, the two-sample analysis with mappability and GC content (Section 5) can be useful. When control sample is not available, 'mosaics' package accommodates one-sample analysis of ChIP-seq data. In this case, you should have files for mappability and GC content, in addition to the files for ChIP and matched control samples.

We recommend users start from Section 3 and it discusses the most convenient way to do the two-sample analysis (without using mappability and GC content). Section 4 discusses each step of the two-sample workflow in detail and provides command lines for each step. Sections 5 and 6 briefly explain the workflow and command lines for the two-sample analysis and the one-sample analysis with mappability and GC content, respectively.

We encourage questions or requests regarding 'mosaics' package to be posted on our Google group [http://groups.google.com/group/mosaics\\_user\\_group](http://groups.google.com/group/mosaics_user_group).

## 3 Two-Sample Analysis using 'mosaicsRunAll'

Two-sample analysis without mappability and GC content can be done in a more convenient way, with the command:

```
R> mosaicsRunAll(  
+   chipFile="/scratch/eland/STAT1_ChIP_eland_results.txt",  
+   chipFileFormat="eland_result",  
+   controlFile="/scratch/eland/STAT1_control_eland_results.txt",  
+   controlFileFormat="eland_result",  
+   binfileDir="/scratch/bin/",  
+   peakFile=c("/scratch/peak/STAT1_peak_list.bed", "/scratch/peak/STAT1_peak_list.gff"),  
+   peakFileFormat=c("bed", "gff"),  
+   reportSummary=TRUE,  
+   summaryFile="/scratch/reports/mosaics_summary.txt",  
+   reportExploratory=TRUE,  
+   exploratoryFile="/scratch/reports/mosaics_exploratory.pdf",  
+   reportGOF=TRUE,  
+   gofFile="/scratch/reports/mosaics_GOF.pdf",  
+   byChr=FALSE,  
+   excludeChr="chrM",  
+   FDR=0.05,  
+   fragLen=200,  
+   binSize=fragLen,  
+   capping=0,  
+   bgEst="automatic",  
+   signalModel="BIC",  
+   parallel=TRUE,  
+   nCore=8 )
```

‘mosaicsRunAll’ method imports aligned read files, converts them to bin-level files (generated bin-level files will be saved in the directory specified in ‘binfileDir’ argument for future use), fits the MOSAiCS model, identifies peaks, and exports the peak list. In addition, users can also make ‘mosaicsRunAll’ method generate diverse analysis reports, such as summary report of parameters and analysis results, exploratory plots, and goodness of fit (GOF) plots. Arguments of ‘mosaicsRunAll’ method are summarized in Table 1. See Section 4.1 for details of the arguments ‘chipFileFormat’, ‘controlFileFormat’, ‘byChr’, ‘excludeChr’, ‘fragLen’, ‘binSize’, and ‘capping’. See Section 4.3 for details of the argument ‘bgEst’. See Section 4.4 for details of the arguments ‘FDR’, ‘signalModel’, ‘peakFileFormat’, ‘maxgap’, ‘minsize’, and ‘thres’.

## 4 Workflow: Two-Sample Analysis

### 4.1 Constructing Bin-Level Files from the Aligned Read File

R package ‘mosaics’ analyzes the data after converting aligned read files into bin-level files for modeling and visualization purposes. These bin-level data can easily be generated from the aligned read files with the command:

```
R> constructBins( infile="/scratch/eland/STAT1_eland_results.txt",
+   fileFormat="eland_result", outfileLoc="/scratch/eland/",
+   byChr=FALSE, excludeChr="chrM", fragLen=200, binSize=200, capping=0 )
```

You can specify the name and file format of the aligned read file in ‘infile’ and ‘fileFormat’ arguments, respectively. ‘constructBins’ method currently allows the following aligned read file formats: Eland result (“eland\_result”), Eland extended (“eland\_extended”), Eland export (“eland\_export”), default Bowtie (“bowtie”), SAM (“sam”), BED (“bed”), and CSEM BED (“csem”). This method assumes that these aligned read files are obtained from single-end tag (SET) experiments. If input file format is neither BED nor CSEM BED, it retains only reads mapping uniquely to the reference genome (uni-reads).

Even though ‘constructBins’ retains only uni-reads for most aligned read file formats, reads mapping to multiple locations on the reference genome (multi-reads) can be easily incorporated into bin-level files by utilizing our multi-read allocator, CSEM (ChIP-Seq multi-read allocator using Expectation-Maximization algorithm). Galaxy tool for CSEM is available in Galaxy Tool Shed (<http://toolshed.g2.bx.psu.edu/>; “csem” under “Next Gen Mappers”). Stand-alone version of CSEM is also available at <http://www.stat.wisc.edu/~keles/Software/multi-reads/>. CSEM exports uni-reads and allocated multi-reads into standard BED file and the corresponding bin-level files can be constructed by applying ‘constructBins’ method to this BED file with the argument ‘fileFormat="csem"’.

‘constructBins’ can generate a single bin-level file containing all chromosomes (for a genome-wide analysis) or multiple bin-level files for each chromosome (for a chromosome-wise analysis). If ‘byChr=FALSE’, bin-level data for all chromosomes are exported to one file named as ‘[infileName]\_fragL[fragLen]\_bin[binSize].txt’, where [infileName], [fragLen], and [binSize] are name of aligned read file, average fragment length, and bin size, respectively. If ‘byChr=TRUE’, bin-level data for each chromosome is exported to a separate file named as ‘[infileName]\_fragL[fragLen]\_bin[binSize]\_chr[chrID].txt’, where [chrID] is chromosome ID that reads align to. These chromosome IDs ([chrID]) are extracted from the aligned read file. The constructed bin-level files are exported to the directory specified in ‘outfileLoc’ argument.

If you want to exclude some chromosomes in the processed bin-level files, you can specify these chromosomes in ‘excludeChr’ argument. You can specify average fragment length and bin size

Table 1: **Summary of the arguments of ‘mosaicsRunAll’ method.**

| (a) Input and output files   |  |
|--|--|
| Argument   | Explanation  |
| <code>chipFile</code>  | Name of aligned read file of ChIP sample.                        |
| <code>chipFileFormat</code>  | File format of aligned read file of ChIP sample.                 |
| <code>controlFile</code>   | Name of aligned read file of matched control sample.             |
| <code>controlFileFormat</code>   | File format of aligned read file of matched control sample.      |
| <code>binfileDir</code>  | Directory that bin-level files are exported to.                  |
| <code>peakFile</code>  | Vector of file names of peak list.                               |
| <code>peakFileFormat</code>  | Vector of file formats of peak list.                             |
| (b) Reports  |  |
| Argument   | Explanation  |
| <code>reportSummary *</code>   | Generate analysis summary?                                       |
| <code>summaryFileName</code>   | File name of analysis summary.                                   |
| <code>reportExploratory *</code>   | Generate exploratory plots?                                      |
| <code>exploratoryFileName</code>   | File name of exploratory plots.                                  |
| <code>reportGOF *</code>   | Generate GOF plots?  |
| <code>gofFileName</code>   | File name of GOF plots.  |
| * Reports will be generated only when these arguments are TRUE. Default is FALSE.                |  |
| (c) Tuning parameters  |  |
| Argument   | Explanation  |
| <code>byChr</code>   | Genome-wide analysis (FALSE) or chromosome-wise analysis (TRUE)? |
| <code>excludeChr</code>  | Vector of chromosomes to be excluded from the analysis.          |
| <code>fragLen</code>   | Average fragment length.   |
| <code>binSize</code>   | Bin size.  |
| <code>capping</code>   | Cap read counts in aligned read files?                           |
| <code>bgEst</code>   | Background estimation approach.                                  |
| <code>signalModel</code>   | Signal model.  |
| <code>FDR</code>   | False discovery rate (FDR).                                      |
| <code>maxgap</code>  | Distance between initial peaks for merging.                      |
| <code>minsize</code>   | Minimum width to be called as a peak.                            |
| <code>thres</code>   | Minimum ChIP tag counts to be called as a peak.                  |
| <code>parallel</code>  | Use <code>parallel</code> package for parallel computing?        |
| <code>nCore</code>   | Number of CPUs used for parallel computing. **                   |
| ** Relevant only when <code>parallel=TRUE</code> and <code>parallel</code> package is installed. |  |

in ‘fragLen’ and ‘binSize’ arguments, respectively, and these arguments control the resolution of bin-level ChIP-seq data. By default, average fragment length is set to 200 bp, which is the common fragment length for Illumina sequences, and bin size equals to average fragment length. ‘capping’ argument indicates maximum number of reads allowed to start at each nucleotide position. Using some small value for capping (e.g., ‘capping=3’) will exclude extremely large read counts that might correspond to PCR amplification artifacts, which is especially useful for the ChIP-seq data with low sequencing depth. Capping is not applied (default) if ‘capping’ is set to some non-positive value, e.g., ‘capping=0’.

## 4.2 Reading Bin-Level Data into the R Environment

You now have bin-level ChIP data and matched control sample data from ‘constructBins’. In this vignette, we use chromosome 21 data from a ChIP-seq experiment of STAT1 binding in interferon- $\gamma$ -stimulated HeLa S3 cells [2]. ‘mosaicsExample’ package provides this example dataset.

```
R> library(mosaicsExample)
```

Bin-level data can be imported to the R environment with the command:

```
R> exampleBinData <- readBins( type=c("chip","input"),
+   fileName=c( system.file( file.path("extdata","chip_chr21.txt"), package="mosaicsExample",
+   system.file( file.path("extdata","input_chr21.txt"), package="mosaicsExample") ) )
```

For the ‘type’ argument, “chip” and “input” indicate bin-level ChIP data control sample data, respectively. You need to specify the corresponding file names in ‘fileName’. ‘mosaics’ package assumes that each file name in ‘fileName’ is provided in the same order as in ‘type’.

In *mosaics* package, you can do either genome-wide analysis or chromosome-wise analysis and this analysis type will be determined automatically based on the contents of bin-level files imported using ‘readBins’. If the bin-level files contain more than one chromosome (i.e., bin-level files are obtained using ‘byChr=FALSE’ in ‘constructBins’), ‘mosaicsFit’ will analyze all the chromosomes simultaneously (genome-wide analysis). Note that if these bin-level files contain different sets of chromosomes, then ‘readBins’ method will utilize only the intersection of them. If bin-level files are obtained using ‘byChr=TRUE’ in ‘constructBins’, each bin-level file contains data for only one chromosome and each of these bin-level files need to be analyzed separately (chromosome-wise analysis). The genome-wide analysis usually provide more stable model fitting and peak identification results so it is recommended for most cases.

R package *mosaics* provides functions for generating simple summaries of the data. The following command prints out basic information about the bin-level data, such as number of bins and total “effective tag counts”. “Total effective tag counts” is defined as the sum of the ChIP tag counts of all bins. This value is usually larger than the sequencing depth since tags are counted after extension to average fragment length and an extended fragment can contribute to multiple bins.

```
R> exampleBinData
```

```
Summary: bin-level data (class: BinData)
```

```
-----
- # of chromosomes in the data: 1
- total effective tag counts: 1637823
```

- (sum of ChIP tag counts of all bins)
- control sample is incorporated
- mappability score is NOT incorporated
- GC content score is NOT incorporated
- uni-reads are assumed

-----

‘print’ method returns the bin-level data in data frame format.

```
R> print(exampleBinData)[51680:51690,]
```

|       | chrID | coord   | tagCount | input |
|-------|-------|---------|----------|-------|
| 51680 | chr21 | 2583950 | 0        | 0     |
| 51681 | chr21 | 2584000 | 0        | 0     |
| 51682 | chr21 | 2584050 | 0        | 0     |
| 51683 | chr21 | 2584100 | 0        | 0     |
| 51684 | chr21 | 2584150 | 0        | 0     |
| 51685 | chr21 | 2584200 | 0        | 0     |
| 51686 | chr21 | 2584250 | 0        | 0     |
| 51687 | chr21 | 2584300 | 0        | 0     |
| 51688 | chr21 | 2584350 | 0        | 0     |
| 51689 | chr21 | 2584400 | 0        | 0     |
| 51690 | chr21 | 2584450 | 0        | 0     |

‘plot’ method provides exploratory plots for the ChIP data. Different type of plots can be obtained by varying the ‘plotType’ argument. ‘plotType="input"’ generates a plot of mean ChIP tag counts versus control tag counts. If ‘plotType’ is not specified, this method plots the histogram of ChIP tag counts.

```
R> plot( exampleBinData )
R> plot( exampleBinData, plotType="input" )
```

Figures 1 and 2 display examples of different types of plots. The relationship between mean ChIP tag counts and control tag counts seems to be linear, especially for small control tag counts (Figure 2).

### 4.3 Fitting the MOSAiCS Model

We are now ready to fit a MOSAiCS model using the bin-level data above (exampleBinData) with the command:

```
R> exampleFit <- mosaicsFit( exampleBinData, analysisType="IO", bgEst="automatic" )
```

‘analysisType="IO"’ indicates implementation of the two-sample analysis. ‘bgEst’ argument determines background estimation approach. ‘bgEst="matchLow"’ estimates background distribution using only bins with low tag counts and it is appropriate for the data with relatively low sequencing depth. ‘bgEst="rMOM"’ estimates background distribution using robust method of moment (MOM) and it is appropriate for the data with relatively high sequencing depth. If ‘bgEst="automatic"’ (default), ‘mosaicsFit’ tries its best guess for the background estimation approach, based on the data provided. If the goodness of fit obtained using ‘bgEst="automatic"’

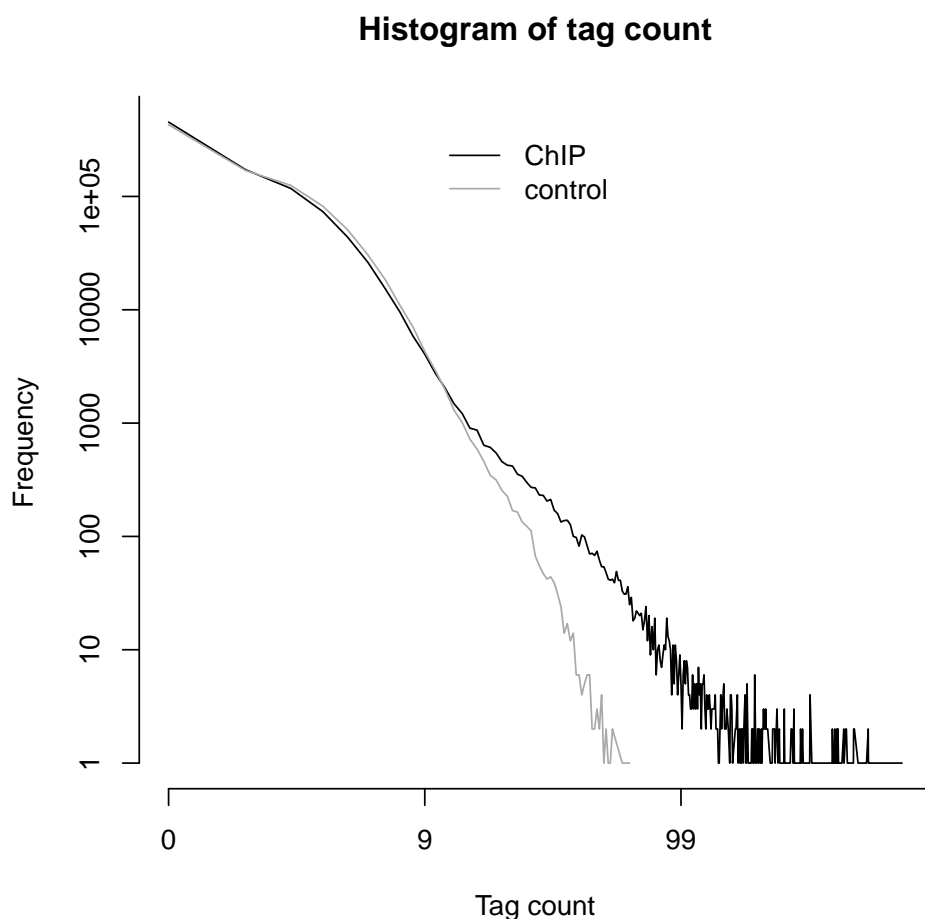


Figure 1: Histograms of the count data from ChIP and control samples.

is not satisfactory, we recommend to try `'bgEst="matchLow"'` and `'bgEst="rMOM"'` and it might improve the model fit.

`'mosaicsFit'` fits both one-signal-component and two-signal-component models. When identifying peaks, you can choose the number of signal components to be used for the final model. The optimal choice of the number of signal components depends on the characteristics of data. In order to support users in the choice of optimal signal model, `mosaics` package provides Bayesian Information Criterion (BIC) values and Goodness of Fit (GOF) plots of these signal models.

The following command prints out BIC values of one-signal-component and two-signal-component models, with additional information about the parameters used in fitting the background (non-enriched) distribution. A lower BIC value indicates a better model fit. For this dataset, we conclude that the two-signal-component model has a lower BIC and hence it provides a better fit.

```
R> exampleFit
```

```
Summary: MOSAiCS model fitting (class: MosaicsFit)
```

```
-----
analysis type: two-sample analysis (Input only)
```

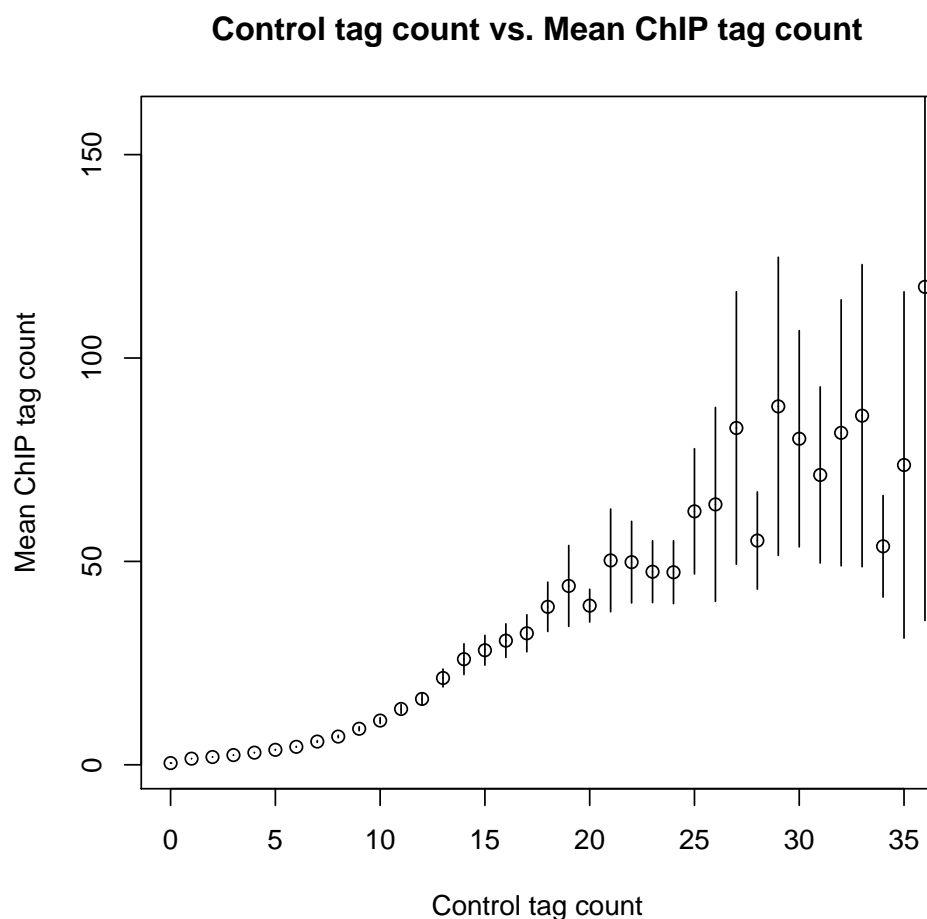


Figure 2: Mean ChIP tag count versus Control tag count.

```
parameters used: k = 3, d = 0.25
BIC of one-signal-component model = 1226948
BIC of two-signal-component model = 1221157
```

-----

‘plot’ method provides the GOF plot. This plots allows visual comparisons of the fits of the background, one-signal-component, and two-signal-component models with the actual data. Figure 3 displays the GOF plot for our dataset and we conclude that the two-signal-component model provides a better fit as is also supported by its lower BIC value compared to the one-signal component model.

```
R> plot(exampleFit)
```

#### 4.4 Identifying Peaks Based on the Fitted Model

Using BIC values and GOF plots in the previous section, we concluded that two-signal-component model fits our data better. Next, we will identify peaks with the two-signal-component model at a



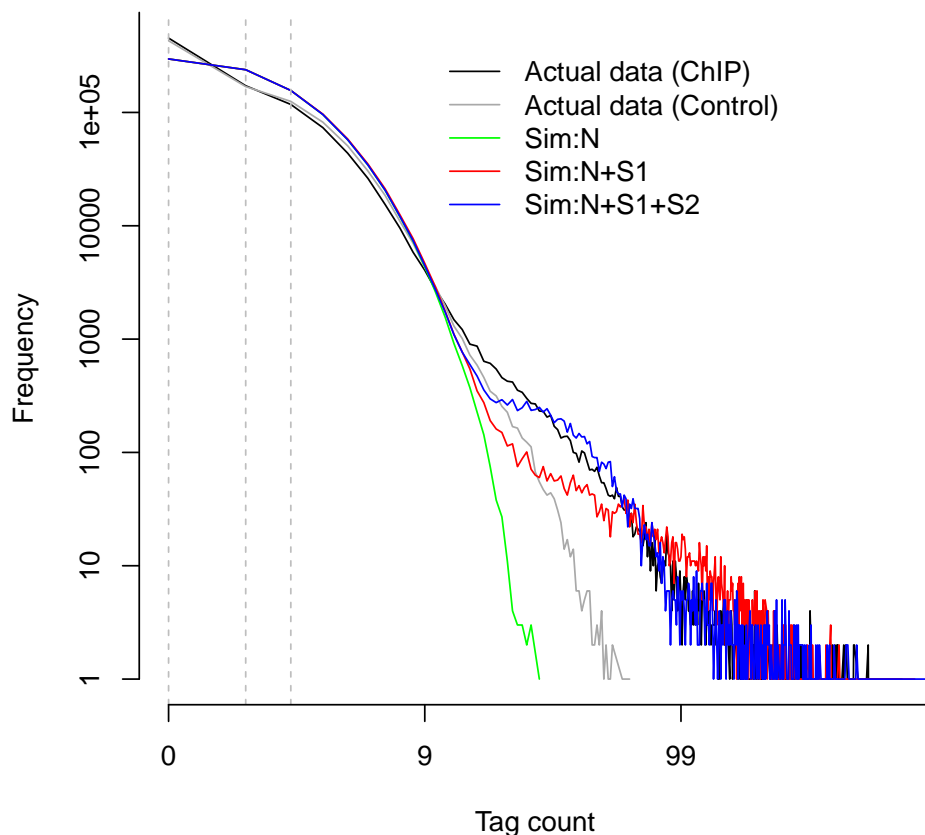


Figure 3: Goodness of Fit (GOF) plot. Depicted are actual data for ChIP and control samples with simulated data from the following fitted models: (Sim:N): Background model; (Sim:N+S1): one-signal-component model; (Sim:N+S1+S2): two-signal-component model.

false discovery rate (FDR) of 0.05 using the command:

```
R> examplePeak <- mosaicsPeak( exampleFit, signalModel="2S", FDR=0.05,
+ maxgap=200, minsize=50, thres=10 )
```

‘signalModel="2S"’ indicates two-signal-component model. Similarly, one-signal-component model can be specified by ‘signalModel="1S"’. FDR can be controlled at the desired level by specifying ‘FDR’ argument. In addition to these two essential parameters, you can also control three more parameters, ‘maxgap’, ‘minsize’, and ‘thres’. These parameters are for refining initial peaks called using specified signal model and FDR. Initial nearby peaks are merged if the distance (in bp) between them is less than ‘maxgap’. Some initial peaks are removed if their lengths are shorter than ‘minsize’ or their ChIP tag counts are less than ‘thres’.

If you use a bin size shorter than the average fragment length in the experiment, we recommend to set ‘maxgap’ to the average fragment length and ‘minsize’ to the bin size. This setting removes peaks that are too narrow (e.g., singletons). If you set the bin size to the average fragment length

(or maybe bin size is larger than the average fragment length), we recommend setting ‘minsize’ to a value smaller than the average fragment length while leaving ‘maxgap’ the same as the average fragment length. This is to prevent filtering using ‘minsize’ because initial peaks would already be at a reasonable width. ‘thres’ is employed to filter out initial peaks with very small ChIP tag counts because such peaks might be false discoveries. Optimal choice of ‘thres’ depends on the sequencing depth of the ChIP-seq data to be analyzed. If you don’t wish to filter out initial peaks using ChIP tag counts, you can set ‘thres’ to an arbitrary negative value.

The following command prints out a summary of identified peaks including the number of peaks identified, median peak width, and the empirical false discovery rate (FDR).

```
R> examplePeak
```

```
Summary: MOSAiCS peak calling (class: MosaicsPeak)
```

```
-----
final model: two-sample analysis (input only) with two signal components
setting: FDR = 0.05, maxgap = 200, minsize = 50, thres = 10
# of peaks = 726
median peak width = 350
empirical FDR = 0.0508
-----
```

‘print’ method returns the peak calling results in data frame format. This data frame can be used as an input for downstream analysis such as motif finding. This output might have different number of columns, depending on ‘analysisType’ of ‘mosaicsFit’. For example, if ‘analysisType="TS"’, columns are peak start position, peak end position, peak width, averaged posterior probability, minimum posterior probability, averaged ChIP tag count, maximum ChIP tag count, averaged control tag count, averaged control tag count scaled by sequencing depth, averaged log base 2 ratio of ChIP over input tag counts, averaged mappability score, and averaged GC content score for each peak. Here, the posterior probability of a bin refers to the probability that the bin is not a peak conditional on data. Hence, smaller posterior probabilities provide more evidence that the bin is actually a peak.

```
R> print(examplePeak)[1:15,]
```

|    | chrID | peakStart | peakStop | peakSize | aveP        | minP         | aveChipCount |
|----|-------|-----------|----------|----------|-------------|--------------|--------------|
| 1  | chr21 | 9874250   | 9874449  | 200      | 0.206469241 | 7.619896e-02 | 25.75000     |
| 2  | chr21 | 10191850  | 10192049 | 200      | 0.199182969 | 6.911445e-02 | 24.75000     |
| 3  | chr21 | 14538150  | 14538499 | 350      | 0.004259010 | 2.726406e-07 | 34.85714     |
| 4  | chr21 | 14677200  | 14677899 | 700      | 0.103764570 | 9.545638e-03 | 29.07143     |
| 5  | chr21 | 14828050  | 14828349 | 300      | 0.030471394 | 1.912252e-03 | 25.66667     |
| 6  | chr21 | 14901600  | 14901799 | 200      | 0.050239923 | 1.406300e-02 | 22.50000     |
| 7  | chr21 | 15175300  | 15175399 | 100      | 0.279052908 | 2.790529e-01 | 19.00000     |
| 8  | chr21 | 15177450  | 15177599 | 150      | 0.171173758 | 6.831582e-02 | 18.66667     |
| 9  | chr21 | 15353150  | 15353549 | 400      | 0.001137400 | 7.910143e-22 | 81.37500     |
| 10 | chr21 | 15358350  | 15358599 | 250      | 0.361641865 | 9.552980e-02 | 18.80000     |
| 11 | chr21 | 15359300  | 15359399 | 100      | 0.180636803 | 1.710742e-01 | 23.50000     |
| 12 | chr21 | 15359950  | 15360349 | 400      | 0.157492156 | 6.213478e-02 | 24.87500     |
| 13 | chr21 | 15374700  | 15375349 | 650      | 0.004989049 | 5.415344e-36 | 93.92308     |
| 14 | chr21 | 15378900  | 15379149 | 250      | 0.068953973 | 2.339854e-02 | 23.00000     |

|    |       |              |               |                     |              |              |          |
|----|-------|--------------|---------------|---------------------|--------------|--------------|----------|
| 15 | chr21 | 15415800     | 15415999      | 200                 | 0.070159735  | 3.788661e-02 | 25.75000 |
|    |       | maxChipCount | aveInputCount | aveInputCountScaled | aveLog2Ratio |              |          |
| 1  |       | 29           | 24.000000     | 26.908521           | -0.04412279  |              |          |
| 2  |       | 29           | 16.750000     | 18.779905           | 0.39916849   |              |          |
| 3  |       | 48           | 2.428571      | 2.722886            | 3.30280262   |              |          |
| 4  |       | 35           | 19.214286     | 21.542834           | 0.44274887   |              |          |
| 5  |       | 31           | 4.666667      | 5.232212            | 2.14088982   |              |          |
| 6  |       | 25           | 3.250000      | 3.643862            | 2.34102405   |              |          |
| 7  |       | 19           | 5.000000      | 5.605942            | 1.59816383   |              |          |
| 8  |       | 20           | 2.333333      | 2.616106            | 2.45618722   |              |          |
| 9  |       | 125          | 4.750000      | 5.325645            | 3.56002727   |              |          |
| 10 |       | 21           | 5.600000      | 6.278655            | 1.45554426   |              |          |
| 11 |       | 24           | 13.000000     | 14.575449           | 0.65695507   |              |          |
| 12 |       | 30           | 14.250000     | 15.976934           | 0.70848160   |              |          |
| 13 |       | 180          | 2.846154      | 3.191075            | 4.22095267   |              |          |
| 14 |       | 25           | 5.800000      | 6.502893            | 1.78233912   |              |          |
| 15 |       | 27           | 11.500000     | 12.893666           | 0.96473565   |              |          |

You can export peak calling results to text files in diverse file formats. Currently, ‘mosaics’ package supports TXT, BED, and GFF file formats. In the exported file, TXT file format (`type="txt"`) includes all the columns that ‘print’ method returns. ‘type="bed"' and ‘type="gff"' export peak calling results in standard BED and GFF file formats, respectively, where score is the averaged ChIP tag counts in each peak. Peak calling results can be exported in TXT, BED, and GFF file formats, respectively, by the commands:

```
R> export( examplePeak, type="txt", filename="TSpeakList.txt" )
R> export( examplePeak, type="bed", filename="TSpeakList.bed" )
R> export( examplePeak, type="gff", filename="TSpeakList.gff" )
```

‘fileLoc’ and ‘fileName’ indicate the directory and the name of the exported file.

## 5 Two-Sample Analysis with Mappability and GC Content

For the two-sample analysis with mappability and GC content and the one-sample analysis, you also need bin-level mappability, GC content, and sequence ambiguity score files for the reference genome you are working with. If you are working with organisms such as human (HG18 and HG19), mouse (MM9), rat (RN4), and Arabidopsis (TAIR9), you can download their corresponding preprocessed mappability, GC content, and sequence ambiguity score files at <http://www.stat.wisc.edu/~keles/Software/mosaics/>. If your reference genome of interest is not listed on our website, you can inquire about it at our Google group, [http://groups.google.com/group/mosaics\\_user\\_group](http://groups.google.com/group/mosaics_user_group), and we would be happy to add your genome of interest to the list. The companion website also provides all the related scripts and easy-to-follow instructions to prepare these files. Please check <http://www.stat.wisc.edu/~keles/Software/mosaics/> for more details.

You can import bin-level data and fit MOSAiCS model for the two-sample analysis using mappability and GC content with the commands:

```
R> exampleBinData <- readBins( type=c("chip","input","M","GC","N"),
+   fileName=c( system.file( file.path("extdata","chip_chr21.txt"), package="mosaicsExample"),
```

```
+ system.file( file.path("extdata","input_chr21.txt"), package="mosaicsExample"),
+ system.file( file.path("extdata","M_chr21.txt"), package="mosaicsExample"),
+ system.file( file.path("extdata","GC_chr21.txt"), package="mosaicsExample"),
+ system.file( file.path("extdata","N_chr21.txt"), package="mosaicsExample") ) )
```

```
-----
Info: preprocessing summary
-----
```

```
- percentage of bins with ambiguous sequences: 27%
  (these bins will be excluded from the analysis)
- before preprocessing:
    first coordinates = 0, last coordinates = 46944350
- after preprocessing:
    first coordinates = 9719550, last coordinates = 46944250
-----
```

For the ‘type’ argument, "chip", "input", "M", "GC", and "N" indicate bin-level ChIP data, control sample data, mappability score, GC content score, and sequence ambiguity score, respectively.

When you have mappability and GC contents, ‘plot’ method provides additional plot types. ‘plotType="M"’ and ‘plotType="GC"’ generate plots of mean ChIP tag counts versus mappability and GC content scores, respectively. Moreover, ‘plotType="M|input”’ and ‘plotType="GC|input”’ generate plots of mean ChIP tag counts versus mappability and GC content scores, respectively, conditional on control tag counts.

```
R> plot( exampleBinData, plotType="M" )
R> plot( exampleBinData, plotType="GC" )
R> plot( exampleBinData, plotType="M|input" )
R> plot( exampleBinData, plotType="GC|input" )
```

As discussed in [1], we observe that mean ChIP tag count increases as mappability score increases (Figure 4). Mean ChIP tag count depends on GC score in a non-linear fashion (Figure 5). When we condition on control tag counts (Figures 6 and 7 ), mean ChIP tag count versus mappability and GC content relations exhibit similar patterns to that of marginal plots given in Figures 4 and 5. MOSAiCS incorporates this observation by modeling ChIP tag counts from non-peak regions with a small number of control tag counts as a function of mappability, GC content, and control tag counts.

Application of MOSAiCS to multiple case studies of ChIP-seq data with low sequencing depth showed that consideration of mappability and GC content in the model improves sensitivity and specificity of peak identification even in the presence of a control sample [1]. `mosaics` package accommodates a two-sample analysis with mappability and GC content by specification of ‘analysisType="TS"’ when calling the ‘`mosaicsFit`’ method.

```
R> exampleFit <- mosaicsFit( exampleBinData, analysisType="TS", bgEst="automatic" )
```

Peak identification can be done exactly in the same way as in the case of the two-sample analysis without mappability and GC content.

```
R> OneSamplePeak <- mosaicsPeak( OneSampleFit, signalModel="2S", FDR=0.05,
+ maxgap=200, minsize=50, thres=10 )
```

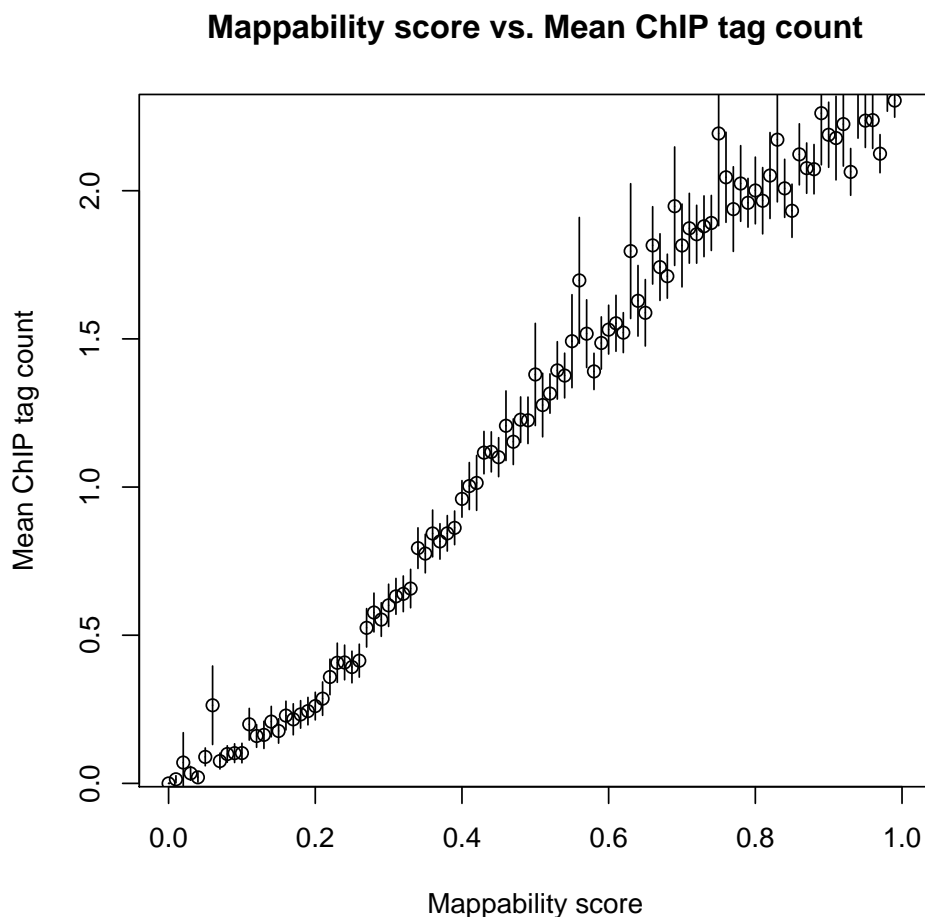


Figure 4: Mean ChIP tag count versus Mappability.

## 6 One-Sample Analysis

When control sample is not available, ‘mosaics’ package accommodates one-sample analysis of ChIP-seq data. Implementation of the MOSAiCS one-sample model is very similar to that of the two-sample analysis. Bin-level data for the one-sample analysis can be imported to the R environment with the command:

```
R> exampleBinData <- readBins( type=c("chip","M","GC","N"),
+   fileName=c( system.file( file.path("extdata","chip_chr21.txt"), package="mosaicsExample"),
+   system.file( file.path("extdata","M_chr21.txt"), package="mosaicsExample"),
+   system.file( file.path("extdata","GC_chr21.txt"), package="mosaicsExample"),
+   system.file( file.path("extdata","N_chr21.txt"), package="mosaicsExample") ) )
```

In order to fit a MOSAiCS model for the one-sample analysis, you need to specify ‘analysisType="OS"’ when calling the ‘mosaicsFit’ method.

```
R> exampleFit <- mosaicsFit( exampleBinData, analysisType="OS", bgEst="automatic" )
```

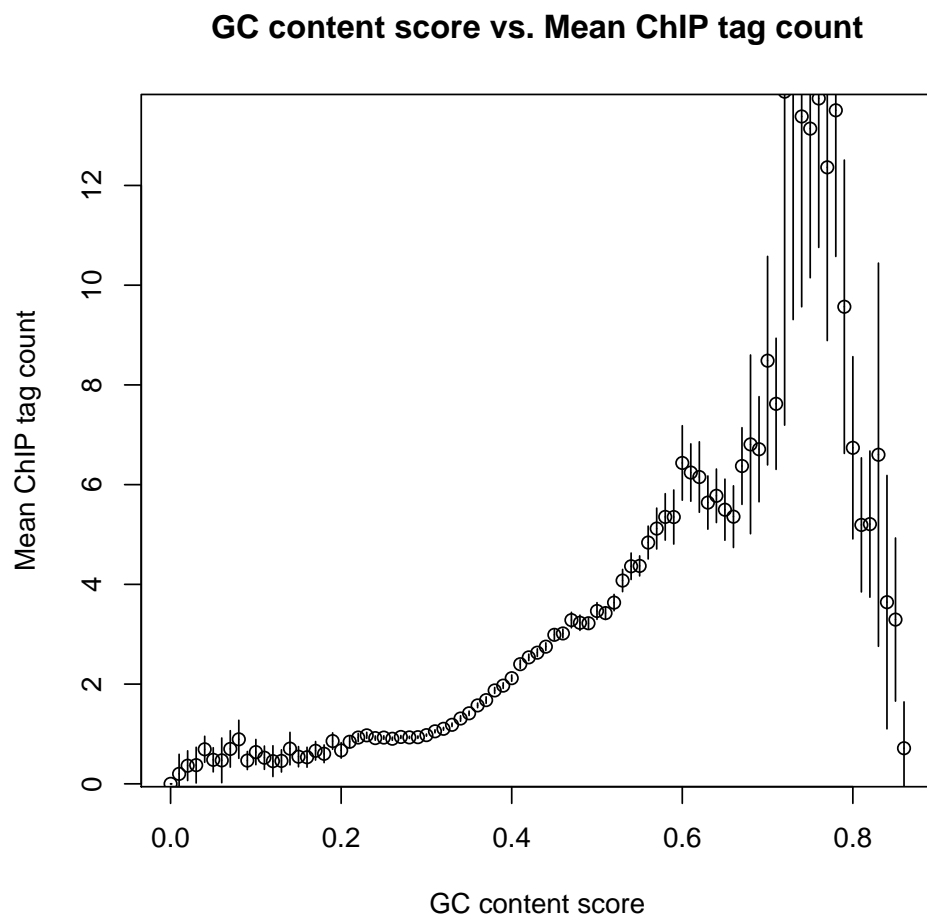


Figure 5: Mean ChIP tag count versus GC content.

Peak identification can be done exactly in the same way as in the case of the two-sample analysis.

```
R> exampleFit <- mosaicsPeak( exampleFit, signalModel="2S", FDR=0.05,
+ maxgap=200, minsize=50, thres=10 )
```

## 7 Conclusion and Ongoing Work

R package `mosaics` provides effective tools to read and investigate ChIP-seq data, fit MOSAiCS model, and identify peaks. We are continuously working on improving `mosaics` package further, especially in supporting more diverse genomes, automating fitting procedures, developing more friendly and easy-to-use user interface, and providing more effective data investigation tools. Please post any questions or requests regarding ‘`mosaics`’ package at [http://groups.google.com/group/mosaics\\_user\\_group](http://groups.google.com/group/mosaics_user_group). Updates and changes of ‘`mosaics`’ package will be announced at our Google group and the companion website (<http://www.stat.wisc.edu/~keles/Software/mosaics/>).

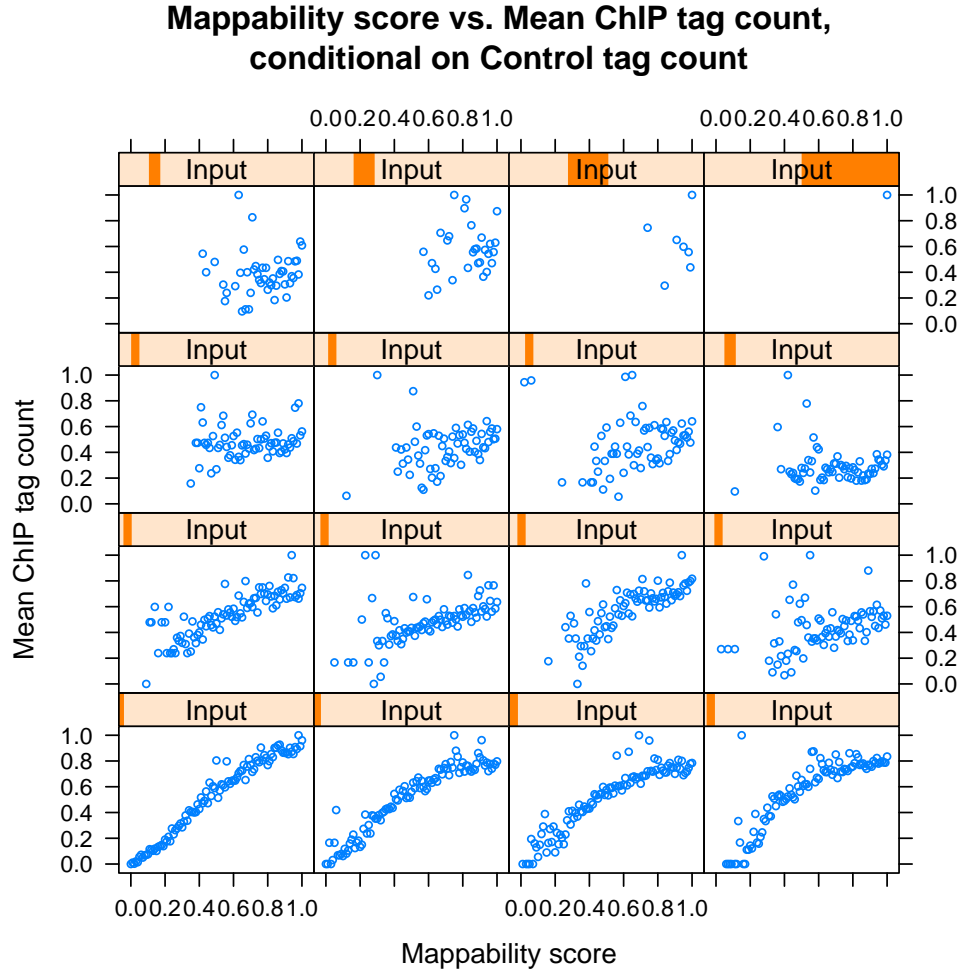


Figure 6: Mean ChIP tag count versus Mappability, conditional on control tag counts.

## References

- [1] Kuan, PF, D Chung, G Pan, JA Thomson, R Stewart, and S Keleş (2010), “A Statistical Framework for the Analysis of ChIP-Seq Data”, *Journal of the American Statistical Association*, 106, 891-903.
- [2] Rozowsky, J, G Euskirchen, R Auerbach, D Zhang, T Gibson, R Bjornson, N Carriero, M Snyder, and M Gerstein (2009), “PeakSeq enables systematic scoring of ChIP-Seq experiments relative to controls”, *Nature Biotechnology*, 27, 66-75.

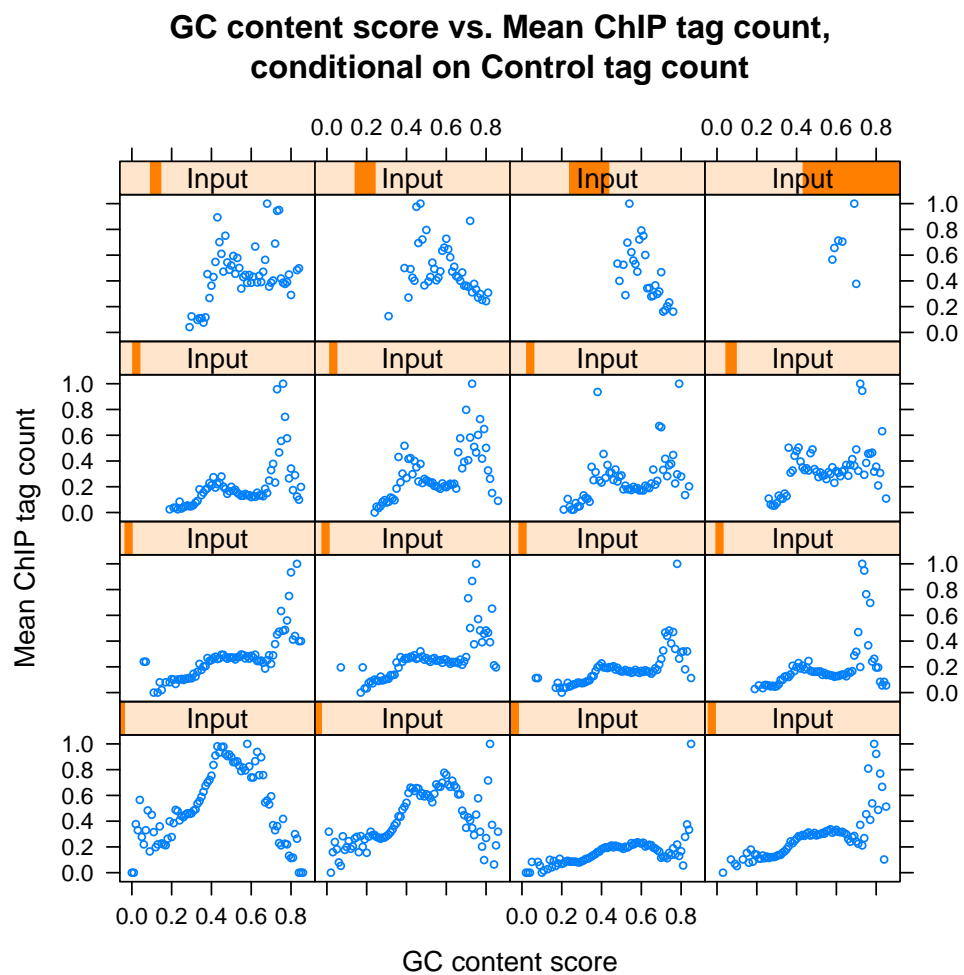


Figure 7: Mean ChIP tag count versus GC content, conditional on control tag counts.