

The ChIPanalyser User's Guide

Patrick Martin

22/08/2017

Introduction

Transcriptional regulation is undeniably a key aspect of cellular homeostasis. It comes to no surprise that modern molecular biology and genomics have showed a keen interest in the subject. Transcription factors (TF) are a force to be reckoned with in the world of transcriptional regulation. Transcription factors are proteins that bind to DNA in a site-specific manner. Experimentally, this binding site can be determined by various methods such as SELEX-seq, EMSA or DNase footprinting. The final result will be a sequence to which a given TF will bind preferentially. In many case, these results are presented in the form of a Position Frequency Matrix or Position Weight Matrix. However at a genome wide scale, modern molecular biology relies on methods such as Chromatin Immuno-precipitation linked to sequencing. This method generates a genome wide profile with peaks at sites of high TF occupancy. These experiments may be very costly and it would be interesting to be able to predict TF occupancy sites *in silico*. With this idea in mind, we present **ChIPanalyser**, a R package developed in the effort of predicting Transcription factor binding. At the core of this package resides an approximation of statistical thermodynamics as suggested by Zabet (Zabet et al. 2015). The statistical thermodynamics framework proposed by Zabet offers a strong ground for binding site prediction as it requires minimal data input. In its current version, ChIPAnalyser requires a DNA sequence, a Position Weight Matrix, the number of bound molecules (or TFs bound to DNA) and a scaling factor for TF specificity. To improve the accuracy of the model, it is also possible to incorporate DNA accessibility data.

Methods

As described above, ChIPAnalyser is based on an approximation of statistical thermodynamics. The core formula describing TF binding is given by :

$$P(N, a, \lambda, \omega)_j = \frac{N \cdot a_j \cdot e^{\left(\frac{1}{\lambda} \cdot \omega_j\right)}}{N \cdot a_j \cdot e^{\left(\frac{1}{\lambda} \cdot \omega_j\right)} + L \cdot n \cdot [a_i \cdot e^{\left(\frac{1}{\lambda} \cdot \omega_j\right)}]_i}$$

with

- N , the number of TF molecules bound to DNA
- a , DNA accessibility
- λ , a parameter scaling the specificity of a given TF
- ω , a Position Weight Matrix.

Work Flow - Quick start

Example data Loading

Before going through the inner workings of the package and the work flow, this section will quickly demonstrate how to load example datasets stored in the package. This data represents a minimal workable examples for the different functions. All data is derived from real biological data in *Drosophila melanogaster* (The *Drosophila melanogaster* genome can be found as a **BSgenome**).

```

library(ChIPAnalyser)

## Loading required package: GenomicRanges
## Loading required package: stats4
## Loading required package: BiocGenerics
## Loading required package: parallel
##
## Attaching package: 'BiocGenerics'
## The following objects are masked from 'package:parallel':
##
##   clusterApply, clusterApplyLB, clusterCall, clusterEvalQ,
##   clusterExport, clusterMap, parApply, parCapply, parLapply,
##   parLapplyLB, parRapply, parSapply, parSapplyLB
## The following objects are masked from 'package:stats':
##
##   IQR, mad, xtabs
## The following objects are masked from 'package:base':
##
##   anyDuplicated, append, as.data.frame, cbind, colnames,
##   do.call, duplicated, eval, evalq, Filter, Find, get, grep,
##   grepl, intersect, is.unsorted, lapply, lengths, Map, mapply,
##   match, mget, order, paste, pmax, pmax.int, pmin, pmin.int,
##   Position, rank, rbind, Reduce, rownames, sapply, setdiff,
##   sort, table, tapply, union, unique, unsplit, which, which.max,
##   which.min
## Loading required package: S4Vectors
##
## Attaching package: 'S4Vectors'
## The following objects are masked from 'package:base':
##
##   colMeans, colSums, expand.grid, rowMeans, rowSums
## Loading required package: IRanges
## Loading required package: GenomeInfoDb
## Loading required package: Biostrings
## Loading required package: XVector
## Loading required package: BSgenome
## Loading required package: rtracklayer
## Loading required package: RcppRoll
#Load data
data(ChIPAnalyserData)

# Loading DNaseSequenceSet from BSgenome object
if(!require("BSgenome.Dmelanogaster.UCSC.dm3", character.only = TRUE)){
  source("https://bioconductor.org/biocLite.R")
}

```

```
biocLite("BSgenome.Dmelanogaster.UCSC.dm3")
}
```

```
## Loading required package: BSgenome.Dmelanogaster.UCSC.dm3
```

```
library(BSgenome.Dmelanogaster.UCSC.dm3)
DNASequenceSet <-getSeq(BSgenome.Dmelanogaster.UCSC.dm3)
```

```
#Loading Position Frequency Matrix
```

```
PFM <- file.path(system.file("extdata",package="ChIPanalyser"),"BCDSLx.pfm")
```

```
#Checking if correctly loaded
ls()
```

```
## [1] "Access"          "DNASequenceSet" "eveLocus"        "eveLocusChip"
## [5] "first_time"      "geneRef"        "PFM"
```

The global environment should now contain a few new variables: DNASequenceSet, PFM, Access, geneRef, eveLocus, eveLocusChip.

- DNASequenceSet is DNASet extracted from the *Drosophila melanogaster* genome (BSgenome). It is advised to use a full genome sequence for this object.
- PFM is a path to file. In this case, it is a Position Frequency Matrix derived from the Bicoid Transcription factor in *Drosophila melanogaster*. This PFM is in RAW format. Although it is possible to directly use a PFM R object, we chose to use a path to a file for this example. Most PFM's downloadable online will come in a text file (with various formats: RAW, TRANSFAC, JASPAR). ChIPanalyser is capable of handling all these formats and parsing these files to usable objects within the package.
- Access is a GRanges object containing accessible DNA for the sequence above.
- geneRef is list of GRanges containing genetic information (exon, intron, 3'UTR, 5'UTR) for the sequence above.
- eveLocus is a GRanges object with genomic position for the eve stripe locus in *Drosophila melanogaster*.
- eveLocusChip is list containing real ChIP-seq data (normalised to each base pair) of the eve stripe locus in *Drosophila melanogaster*.

This section presents a quick work flow. For details on the work flow and objects, see section **Work Flow - Full Guide**

Quick Start

Step 1 - Building Data objects

The first step is to set up your data storing objects. These objects will automatically compute Position Weight Matrix from a Position Frequency Matrix, and Base Pair Frequency from a DNASet. The values that are provided in this example are extracted from real biological data.

NOTE: These values will differ depending on the source of the data and the data itself.

```
# Building a genomicProfileParameters objects for data
# storage and PWM computation
GPP <- genomicProfileParameters(PFM=PFM,PFMFormat="raw",
  BPFrequency=DNASequenceSet,
  ScalingFactorPWM = 1.5,
```

```
PWMThreshold = 0.7)
GPP
```

```
## Object Class:genomicProfileParameters
##
##
## PWM:
##      [,1]      [,2]      [,3]      [,4]      [,5]      [,6]      [,7]
## A  0.1267378 -0.8713677 -3.953162  1.983869  1.983869 -5.052697 -9.445015
## C  0.2913871  0.6224195 -4.801159 -9.445015 -9.445015 -9.445015  1.998447
## G  0.3703684 -2.3054635 -9.445015 -4.587034 -4.587034 -3.422647 -9.445015
## T -1.3522577  0.7753635  1.962784 -9.445015 -9.445015  1.954263 -9.445015
##      [,8]
## A -4.235561
## C  1.831691
## G -3.830305
## T -1.657112
##
## PFM:
##      [,1] [,2] [,3] [,4] [,5] [,6] [,7] [,8]
## A  190   95   11  689  689    5    0    9
## C  213  268    6    0    0    0  696  620
## G  225   35    0    7    7   16    0   12
## T   68  298  679    0    0  675    0   55
##
## PFMFormat: raw
##
## PWM Scores at Sites higher than Threshold:
## GRangesList object of length 0:
## <0 elements>
##
## -----
## seqinfo: no sequences
##
## No Accessible DNA at Loci:
##
## Genomic Profile Parameters:
## Lambda: 1.5
## BP Frequency: 0.25 0.25 0.25 0.25
## Pseudocount: 1
## Natural log: FALSE
## Number Of Sites: 0
## maxPWMScore:
## minPWMScore:
## PWMThreshold: 0.7
## Average Exponential PWM Score:
## DNA Sequence Length:
```

```
## Strand Rule: max
## Strand: +-
# Building occupancyProfileParameters with default values
OPP <- occupancyProfileParameters()
OPP

## Object Class:occupancyProfileParameters
##

## Ploidy: 2
## boundMolecules: 1000
## backgroundSignal: 0
## maxSignal: 1
## chipMean: 150
## chipSd: 150
## chipSmooth: 250
## Step Size: 10
## Theta Threshold: 0.1
# Building occupancyProfileParameters with custom values
OPP <- occupancyProfileParameters(ploidy= 2,
  boundMolecules= 1000,
  chipMean = 200,
  chipSd = 200,
  chipSmooth = 250,
  maxSignal = 1.847,
  backgroundSignal = 0.02550997)
OPP

## Object Class:occupancyProfileParameters
##

## Ploidy: 2
## boundMolecules: 1000
## backgroundSignal: 0.02550997
## maxSignal: 1.847
## chipMean: 200
## chipSd: 200
## chipSmooth: 250
## Step Size: 10
## Theta Threshold: 0.1
```

Step 2 - Optimal Parameters

The model is based on the approximation of statistical thermodynamics with inference of two parameters (ScalingFactorPWM and boundMolecules). In order to infer these parameters, we suggest to use `computeOptimal`. Values that should be tested for `ScalingFactorPWM` and for `boundMolecules` should be provided by user as described above. If these values are not provided (default value OR only one value for each parameter), then they will be assigned internally. The internal values are the following:

```
ScalingFactorPWM(genomicProfileParameters) <- c(0.25, 0.5, 0.75, 1, 1.25,
  1.5, 1.75, 2, 2.5, 3, 3.5 ,4 ,4.5, 5)

boundMolecules(occupancyProfileParameters) <- c(1, 10, 20, 50, 100,
```

```
200, 500,1000,2000, 5000,10000,20000,50000, 100000,
200000, 500000, 1000000)
```

computeOptimalcontains the following arguments:

```
optimalParam <- computeOptimal(DNASequenceSet = DNASequenceSet,
  genomicProfileParameters = GPP,
  LocusProfile = eveLocusChip,
  setSequence = eveLocus,
  DNAAccessibility = Access,
  occupancyProfileParameters = OPP,
  parameter = "all",
  peakMethod="moving_kernel")
```

```
## Computing Genome Wide PWM Score
```

```
## Computing PWM Score at Loci & Extracting Sites Above Threshold
```

```
## Computing Occupancy
```

```
## Computing ChIP-seq-like Profile
```

```
## Computing Accuracy of Profile
```

```
## Extracting Optimal Set of Parameters
```

```
optimalParam
```

```
## $`Optimal Parameters`
```

```
## $`Optimal Parameters`$meanCorr
```

```
## [1] "0.75" "5e+05"
```

```
##
```

```
## $`Optimal Parameters`$meanMSE
```

```
## [1] "1.25" "1000"
```

```
##
```

```
## $`Optimal Parameters`$meanTheta
```

```
## [1] "1.25" "1000"
```

```
##
```

```
##
```

```
## $`Optimal Matrix`
```

```
## $`Optimal Matrix`$meanCorr
```

```
##          1          10          20          50          100          200          500
## 0.25 0.8435269 0.8342830 0.8093797 0.7578630 0.7256334 0.7109625 0.7143568
## 0.5  0.8242734 0.8191801 0.8076173 0.7831160 0.7663309 0.7572188 0.7576516
## 0.75 0.8429730 0.8444026 0.8392846 0.8277879 0.8123382 0.8003016 0.8010699
## 1    0.7379446 0.7689522 0.7885026 0.8099659 0.8151623 0.8142855 0.8164442
## 1.25 0.7671239 0.7900837 0.8064706 0.8280682 0.8352377 0.8380781 0.8345198
## 1.5  0.7283106 0.7496309 0.7725609 0.8026694 0.8226631 0.8277529 0.8357009
## 1.75 0.5823179 0.6035664 0.6263213 0.6781727 0.7290164 0.7711215 0.8038425
## 2    0.4845802 0.4948444 0.5119707 0.5612225 0.6138851 0.6853372 0.7627732
## 2.5  0.3647584 0.3726325 0.3772826 0.3992819 0.4491398 0.5201538 0.6034785
## 3    0.3005580 0.3063813 0.3127481 0.3312094 0.3599365 0.4104222 0.5189434
## 3.5  0.3003317 0.3041450 0.3083337 0.3206006 0.3400799 0.3756586 0.4599116
## 4    0.3001897 0.3027362 0.3055418 0.3138110 0.3271125 0.3520083 0.4148106
## 4.5  0.3001006 0.3018503 0.3037820 0.3094991 0.3187735 0.3364101 0.3827839
## 5    0.3000437 0.3012842 0.3026555 0.3067258 0.3133666 0.3261313 0.3606474
##          1000          2000          5000          10000          20000          50000          1e+05
## 0.25 0.7236881 0.7311739 0.7364203 0.7405083 0.7483979 0.7683033 0.7896878
```

```

## 0.5 0.7652147 0.7777754 0.8004981 0.8175934 0.8304955 0.8377855 0.8400819
## 0.75 0.8061322 0.8145761 0.8287463 0.8354265 0.8384978 0.8402069 0.8415842
## 1 0.8203381 0.8287708 0.8379507 0.8420426 0.8445218 0.8461446 0.8447290
## 1.25 0.8367129 0.8387343 0.8419233 0.8443199 0.8445630 0.8375789 0.8246078
## 1.5 0.8368432 0.8377148 0.8390813 0.8384984 0.8348416 0.8192797 0.7993414
## 1.75 0.8161703 0.8278558 0.8318925 0.8306914 0.8234392 0.8002894 0.7727476
## 2 0.7964770 0.8138327 0.8243079 0.8208994 0.8070174 0.7730552 0.7427729
## 2.5 0.6798135 0.7436438 0.7821902 0.7869273 0.7746124 0.7349417 0.6900998
## 3 0.6168327 0.6966421 0.7516581 0.7609388 0.7481344 0.7015738 0.6478850
## 3.5 0.5494787 0.6376305 0.7116756 0.7302299 0.7218538 0.6755899 0.6194948
## 4 0.4897657 0.5760687 0.6649005 0.6951890 0.6949361 0.6547311 0.6006854
## 4.5 0.4426567 0.5202375 0.6157053 0.6571144 0.6667172 0.6361429 0.5869749
## 5 0.4076712 0.4740565 0.5685194 0.6181697 0.6374511 0.6180478 0.5754023
## 2e+05 5e+05 1e+06
## 0.25 0.8109340 0.8295861 0.8358922
## 0.5 0.8400572 0.8384581 0.8369426
## 0.75 0.8432930 0.8463342 0.8440491
## 1 0.8385323 0.8182547 0.7971591
## 1.25 0.8057857 0.7798274 0.7642323
## 1.5 0.7777255 0.7467306 0.7165258
## 1.75 0.7405593 0.6941079 0.6454311
## 2 0.7049411 0.6382310 0.5761837
## 2.5 0.6334533 0.5456158 0.4799302
## 3 0.5827595 0.4917561 0.4309691
## 3.5 0.5525586 0.4635016 0.4072382
## 4 0.5352262 0.4490886 0.3958365
## 4.5 0.5248570 0.4418802 0.3905427
## 5 0.5178524 0.4383425 0.3883538
##
## $`Optimal Matrix`$meanMSE
## 1 10 20 50 100 200 500
## 0.25 15.12395 14.21461 13.23944 10.79892 8.246195 6.338610 7.188663
## 0.5 15.14775 14.49972 13.80666 11.97720 9.698463 7.069336 5.224428
## 0.75 15.06633 14.49815 13.90746 12.30614 10.266348 7.609517 4.694015
## 1 14.87247 14.32789 13.78662 12.40288 10.588732 8.085813 4.852495
## 1.25 14.94691 14.48733 14.01461 12.79316 11.160155 8.729014 5.208727
## 1.5 14.85102 14.44607 14.02404 12.95486 11.513258 9.380928 5.830681
## 1.75 14.82695 14.44046 14.04332 13.11677 11.874841 9.936094 6.549153
## 2 14.73011 14.41764 14.09193 13.29056 12.176645 10.502878 7.362967
## 2.5 14.84301 14.66748 14.47833 13.90115 12.938717 11.612982 9.170871
## 3 14.88146 14.78166 14.67103 14.34096 13.798384 12.751281 10.673814
## 3.5 14.88582 14.82522 14.75792 14.55638 14.222098 13.562704 11.845118
## 4 14.88832 14.85014 14.80773 14.68053 14.468804 14.047445 12.818200
## 4.5 14.88977 14.86469 14.83682 14.75320 14.613830 14.335498 13.510465
## 5 14.89065 14.87346 14.85436 14.79704 14.701483 14.510387 13.940021
## 1000 2000 5000 10000 20000 50000 1e+05
## 0.25 9.898730 12.722131 15.303409 16.160629 16.068934 14.613783 12.861865
## 0.5 6.127293 7.909539 9.215280 9.201580 8.947431 8.763991 8.737897
## 0.75 4.460814 5.588840 7.166123 7.894866 8.393396 8.835339 9.170573
## 1 3.959862 4.631050 6.400719 7.546960 8.488909 9.765428 11.453707
## 1.25 3.692721 3.997591 6.026234 7.671015 9.514157 13.214493 17.684895
## 1.5 3.835814 3.722497 6.022811 8.607104 11.806207 18.343158 25.076909
## 1.75 4.342797 3.768624 6.216138 9.611625 14.219735 23.471404 33.121121
## 2 4.988651 3.984794 6.064756 10.436707 17.246568 30.303889 41.877836

```

```

## 2.5 7.017935 5.471651 6.941999 11.963012 20.689945 38.389298 56.772459
## 3 8.629886 6.522711 6.147709 10.394305 20.173546 42.593175 67.033802
## 3.5 10.135963 7.997054 6.142466 8.525504 17.386348 41.887097 70.376460
## 4 11.335484 9.438444 6.878056 7.334774 13.844806 37.418655 67.783001
## 4.5 12.219774 10.638032 7.964931 7.071300 10.878793 31.155742 61.290223
## 5 13.014314 11.564158 9.078687 7.470439 9.019120 24.824402 52.902116
## 2e+05 5e+05 1e+06
## 0.25 11.159997 9.688598 9.155631
## 0.5 8.773874 8.954094 9.121116
## 0.75 9.630268 10.814928 12.668697
## 1 14.212007 20.329057 26.240616
## 1.25 23.596167 32.031058 37.280511
## 1.5 32.317903 43.134382 54.138043
## 1.75 44.458670 61.447404 80.364537
## 2 56.070262 82.281302 109.809345
## 2.5 80.609722 122.491255 160.148458
## 3 98.379356 149.286541 167.898645
## 3.5 106.705068 162.728866 163.271534
## 4 107.280293 163.918908 160.254837
## 4.5 102.483718 157.399812 157.278501
## 5 94.352806 150.126047 153.791510
##
## $`Optimal Matrix`$meanTheta
## 1 10 20 50 100 200 500
## 0.25 0.1559489 0.1542399 0.1496359 0.1463682 0.1463682 0.1463682 0.1463682
## 0.5 0.1523894 0.1514477 0.1493100 0.1463682 0.1463682 0.1463682 0.1515391
## 0.75 0.1558465 0.1561108 0.1551646 0.1530391 0.1501828 0.1479575 0.1706577
## 1 0.1463682 0.1463682 0.1463682 0.1497442 0.1507049 0.1505428 0.1682525
## 1.25 0.1463682 0.1463682 0.1490980 0.1530910 0.1544164 0.1549415 0.1602157
## 1.5 0.1463682 0.1463682 0.1463682 0.1483953 0.1520917 0.1530327 0.1545020
## 1.75 0.1463682 0.1463682 0.1463682 0.1463682 0.1463682 0.1463682 0.1486122
## 2 0.1463682 0.1463682 0.1463682 0.1463682 0.1463682 0.1463682 0.1463682
## 2.5 0.1463682 0.1463682 0.1463682 0.1463682 0.1463682 0.1463682 0.1463682
## 3 0.1463682 0.1463682 0.1463682 0.1463682 0.1463682 0.1463682 0.1463682
## 3.5 0.1463682 0.1463682 0.1463682 0.1463682 0.1463682 0.1463682 0.1463682
## 4 0.1463682 0.1463682 0.1463682 0.1463682 0.1463682 0.1463682 0.1463682
## 4.5 0.1463682 0.1463682 0.1463682 0.1463682 0.1463682 0.1463682 0.1463682
## 5 0.1463682 0.1463682 0.1463682 0.1463682 0.1463682 0.1463682 0.1463682
## 1000 2000 5000 10000 20000 50000 1e+05
## 0.25 0.1463682 0.1463682 0.1463682 0.1463682 0.1463682 0.1463682 0.1463682
## 0.5 0.1463682 0.1463682 0.1479939 0.1511544 0.1535397 0.1548875 0.1553120
## 0.75 0.1807142 0.1505966 0.1532163 0.1544513 0.1550191 0.1553351 0.1555897
## 1 0.2071633 0.1789596 0.1549180 0.1556745 0.1561328 0.1564329 0.1561712
## 1.25 0.2265844 0.2098099 0.1556524 0.1560955 0.1561405 0.1548493 0.1524512
## 1.5 0.2181657 0.2250411 0.1551270 0.1550192 0.1543432 0.1514662 0.1477800
## 1.75 0.1879365 0.2196705 0.1537980 0.1535759 0.1522351 0.1479553 0.1463682
## 2 0.1596578 0.2042346 0.1523957 0.1517656 0.1491991 0.1463682 0.1463682
## 2.5 0.1463682 0.1463682 0.1463682 0.1463682 0.1463682 0.1463682 0.1463682
## 3 0.1463682 0.1463682 0.1463682 0.1463682 0.1463682 0.1463682 0.1463682
## 3.5 0.1463682 0.1463682 0.1463682 0.1463682 0.1463682 0.1463682 0.1463682
## 4 0.1463682 0.1463682 0.1463682 0.1463682 0.1463682 0.1463682 0.1463682
## 4.5 0.1463682 0.1463682 0.1463682 0.1463682 0.1463682 0.1463682 0.1463682
## 5 0.1463682 0.1463682 0.1463682 0.1463682 0.1463682 0.1463682 0.1463682
## 2e+05 5e+05 1e+06

```



```
## 0.25 0.1499232 0.1533716 0.1545374
## 0.5 0.1553074 0.1550118 0.1547316
## 0.75 0.1559057 0.1564679 0.1560455
## 1 0.1550255 0.1512767 0.1473766
## 1.25 0.1489714 0.1463682 0.1463682
## 1.5 0.1463682 0.1463682 0.1463682
## 1.75 0.1463682 0.1463682 0.1463682
## 2 0.1463682 0.1463682 0.1463682
## 2.5 0.1463682 0.1463682 0.1463682
## 3 0.1463682 0.1463682 0.1463682
## 3.5 0.1463682 0.1463682 0.1463682
## 4 0.1463682 0.1463682 0.1463682
## 4.5 0.1463682 0.1463682 0.1463682
## 5 0.1463682 0.1463682 0.1463682
##
##
## $Parameter
## [1] "all"
```

This Function might take some time to compute. Do not be alarmed if it takes some time to run. You should be notified of the progress of the function as it goes

This function is a combination of all the functions bellow with some more magic to it. In the following steps we will describe each of the functions.

Step 3 - Genome Wide Scoring

Computing Genome Wide metrics that will be used further down the line.

```
genomeWide <- computeGenomeWidePWMScore(DNASequenceSet=DNASequenceSet,
    genomicProfileParameters=GPP, DNAAccessibility = Access)
```

```
## Scoring whole genome
```

```
## Accessible DNA ~ Both strands
```

```
## Computing Mean waiting time
```

```
genomeWide
```

```
## Object Class:genomicProfileParameters
```

```
##
```

```
##
```

```
## PWM:
```

```
##      [,1]      [,2]      [,3]      [,4]      [,5]      [,6]      [,7]
## A  0.1267378 -0.8713677 -3.953162  1.983869  1.983869 -5.052697 -9.445015
## C  0.2913871  0.6224195 -4.801159 -9.445015 -9.445015 -9.445015  1.998447
## G  0.3703684 -2.3054635 -9.445015 -4.587034 -4.587034 -3.422647 -9.445015
## T -1.3522577  0.7753635  1.962784 -9.445015 -9.445015  1.954263 -9.445015
##      [,8]
## A -4.235561
## C  1.831691
## G -3.830305
## T -1.657112
##
```

```
## PFM:
##      [,1] [,2] [,3] [,4] [,5] [,6] [,7] [,8]
## A   190   95   11  689  689    5    0    9
## C   213  268    6    0    0    0  696  620
## G   225   35    0    7    7   16    0   12
## T    68  298  679    0    0  675    0   55
##
## PFMFormat: raw
##
## PWM Scores at Sites higher than Threshold:
## GRangesList object of length 0:
## <0 elements>
##
## -----
## seqinfo: no sequences
##
## No Accessible DNA at Loci:
##
## Genomic Profile Parameters:
## Lambda: 1.5
## BP Frequency:    0.25    0.25    0.25    0.25
## Pseudocount: 1
## Natural log: FALSE
## Number Of Sites: 0
## maxPWMScore: 12.8606543674325
## minPWMScore: -48.8262800777777
## PWMThreshold: 0.7
## Average Exponential PWM Score: 1.015637
## DNA Sequence Length: 3112514
## Strand Rule: max
## Strand: +-
computeGenomeWidePWMScore will return a genomicProfileParameters object with updated values for
maxPWMScore, minPWMScore, averageExpPWMScore, and DNASequencLength.
```

Step 4 - PWM Scores Above Threshold

Once genome wide scores have been computed, the `genomeWide` object (previously computed) should be parsed to the next function. The next function will compute sites above the assigned threshold (see below) for a given locus (or set of loci). If no Locus is provided then the whole genome will be considered.

```
SitesAboveThreshold <- computePWMScore(DNASequencSet=DNASequencSet,
  genomicProfileParameters=genomeWide,
  setSequence=eveLocus, DNAAccessibility = Access)
```

```
## Processing DNA Acccssibility
## Extracting Sites Above threshold
SitesAboveThreshold
```

```

## Object Class:genomicProfileParameters
##

##
## PWM:

##      [,1]      [,2]      [,3]      [,4]      [,5]      [,6]      [,7]
## A  0.1267378 -0.8713677 -3.953162  1.983869  1.983869 -5.052697 -9.445015
## C  0.2913871  0.6224195 -4.801159 -9.445015 -9.445015 -9.445015  1.998447
## G  0.3703684 -2.3054635 -9.445015 -4.587034 -4.587034 -3.422647 -9.445015
## T -1.3522577  0.7753635  1.962784 -9.445015 -9.445015  1.954263 -9.445015
##      [,8]
## A -4.235561
## C  1.831691
## G -3.830305
## T -1.657112

##
## PFM:

##      [,1] [,2] [,3] [,4] [,5] [,6] [,7] [,8]
## A   190   95   11  689  689    5    0    9
## C   213  268    6    0    0    0  696  620
## G   225   35    0    7    7   16    0   12
## T    68  298  679    0    0  675    0   55

##
## PFMFormat: raw

##
## PWM Scores at Sites higher than Threshold:

## GRangesList object of length 1:
## $eve
## GRanges object with 412 ranges and 1 metadata column:
##      seqnames      ranges strand |      PWMScore
##      <Rle>      <IRanges> <Rle> |      <numeric>
##      [1]    chr2R [5860705, 5860712]    + | -1.51655573585429
##      [2]    chr2R [5860709, 5860716]    + | -5.33217184502491
##      [3]    chr2R [5860715, 5860722]    + |  9.13992557549757
##      [4]    chr2R [5860728, 5860735]    + |  5.05434682102833
##      [5]    chr2R [5860758, 5860765]    + | -5.15370980167748
##      ...      ...      ...      ...      ...
##      [408] chr2R [5876629, 5876636]    + |  5.60817413411963
##      [409] chr2R [5876635, 5876642]    + |  0.202790199774102
##      [410] chr2R [5876641, 5876648]    - | -4.47385601266488
##      [411] chr2R [5876666, 5876673]    + |  2.21133362723558
##      [412] chr2R [5876684, 5876691]    + | -2.28895797651261
##
## -----
## seqinfo: 1 sequence from an unspecified genome; no seqlengths
##
## No Accessible DNA at Loci:
## -
##
## Genomic Profile Parameters:

```

```
## Lambda: 1.5
## BP Frequency: 0.25 0.25 0.25 0.25

## Pseudocount: 1
## Natural log: FALSE
## Number Of Sites: 0
## maxPWMScore: 12.8606543674325
## minPWMScore: -48.8262800777777
## PWMThreshold: 0.7

## Average Exponential PWM Score: 1.015637

## DNA Sequence Length: 3112514
## Strand Rule: max
## Strand: +-

```

This function returns another `genomicProfileParameters` object with an updated `AllSitesAboveThreshold` slot. This slot contains a `GRanges` object with sites above threshold and associated PWMScores.

Step 4 - compute Occupancy

From the PWMScores, ChIPanalyser will compute occupancy for each sites above threshold.

```
Occupancy <- computeOccupancy(SitesAboveThreshold,
  occupancyProfileParameters= OPP)

```

```
## Computing Occupancy at sites higher than threshold.

```

```
Occupancy

```

```
## Object Class:genomicProfileParameters
##

##
## PWM:

##      [,1]      [,2]      [,3]      [,4]      [,5]      [,6]      [,7]
## A  0.1267378 -0.8713677 -3.953162  1.983869  1.983869 -5.052697 -9.445015
## C  0.2913871  0.6224195 -4.801159 -9.445015 -9.445015 -9.445015  1.998447
## G  0.3703684 -2.3054635 -9.445015 -4.587034 -4.587034 -3.422647 -9.445015
## T -1.3522577  0.7753635  1.962784 -9.445015 -9.445015  1.954263 -9.445015
##      [,8]
## A -4.235561
## C  1.831691
## G -3.830305
## T -1.657112

##
## PFM:

##      [,1] [,2] [,3] [,4] [,5] [,6] [,7] [,8]
## A   190   95   11  689  689    5    0    9
## C   213  268    6    0    0    0  696  620
## G   225   35    0    7    7   16    0   12
## T    68  298  679    0    0  675    0   55

##
## PFMFormat: raw

```

```

##
## PWM Scores at Sites higher than Threshold:
## $`lambda = 1.5 & boundMolecules = 1000`
## GRangesList object of length 1:
## $eve
## GRanges object with 412 ranges and 2 metadata columns:
##      seqnames      ranges strand |      PWMScore
##      <Rle>      <IRanges> <Rle> |      <numeric>
## eve chr2R [5860705, 5860712] + | -1.51655573585429
## eve chr2R [5860709, 5860716] + | -5.33217184502491
## eve chr2R [5860715, 5860722] + |  9.13992557549757
## eve chr2R [5860728, 5860735] + |  5.05434682102833
## eve chr2R [5860758, 5860765] + | -5.15370980167748
## ...      ...      ...      ...      ...
## eve chr2R [5876629, 5876636] + |  5.60817413411963
## eve chr2R [5876635, 5876642] + |  0.202790199774102
## eve chr2R [5876641, 5876648] - | -4.47385601266488
## eve chr2R [5876666, 5876673] + |  2.21133362723558
## eve chr2R [5876684, 5876691] + | -2.28895797651261
##      Occupancy
##      <numeric>
## eve 0.0138683203024566
## eve 0.0138160293072631
## eve 0.0783704718441574
## eve 0.0183246335750422
## eve 0.0138165926852252
## ...      ...
## eve 0.0203268664876583
## eve 0.0139901018008407
## eve 0.0138194725807681
## eve 0.014492381648295
## eve 0.0138454813805688
##
## -----
## seqinfo: 1 sequence from an unspecified genome; no seqlengths
##
## No Accessible DNA at Loci:
## -
##
## Genomic Profile Parameters:
## Lambda: 1.5
## BP Frequency: 0.25 0.25 0.25 0.25
## Pseudocount: 1
## Natural log: FALSE
## Number Of Sites: 0
## maxPWMScore: 12.8606543674325
## minPWMScore: -48.8262800777777
## PWMThreshold: 0.7
## Average Exponential PWM Score: 1.015637
## DNA Sequence Length: 3112514

```

```
## Strand Rule: max
## Strand: +-
```

This function will return a `genomicProfileParameters` object with an updated `AllSitesAboveThreshold`. Now the Occupancy values for each sites are included.

Step 5 - compute ChIP -seq like profiles

The ultimate goal of `ChIPAnalyser` is to produce ChIP-seq like profile predicting transcription factor binding. To do so, the following function will compute ChIP-seq like scores from occupancy values.

```
chipProfile <- computeChipProfile(setSequence = eveLocus,
  occupancy = Occupancy, occupancyProfileParameters = OPP,
  method="moving_kernel")
```

```
## Computing ChIP Profile
```

```
chipProfile
```

```
## $`lambda` = 1.5 & boundMolecules = 1000`
## $`lambda` = 1.5 & boundMolecules = 1000`$eve
## GRanges object with 1600 ranges and 1 metadata column:
##      seqnames      ranges strand |      ChIP
##      <Rle>      <IRanges> <Rle> |      <numeric>
## eve chr2R [5860693, 5860703] * | 0.0467998729244692
## eve chr2R [5860703, 5860713] * | 0.051053053031132
## eve chr2R [5860713, 5860723] * | 0.0554324704104929
## eve chr2R [5860723, 5860733] * | 0.059949075887137
## eve chr2R [5860733, 5860743] * | 0.0646141633273505
## ...      ...      ...      ... | ...
## eve chr2R [5876643, 5876653] * | 0.0158243321666681
## eve chr2R [5876653, 5876663] * | 0.0149496630784295
## eve chr2R [5876663, 5876673] * | 0.0140710281760898
## eve chr2R [5876673, 5876683] * | 0.0131862304147329
## eve chr2R [5876683, 5876693] * | 0.0122930573390848
## -----
## seqinfo: 1 sequence from an unspecified genome; no seqlengths
```

This function will return a `List` of `GRangesLists` of `GRanges`. Each element of the list represents a combination of `ScalingFactorPWM` and `boundMolecules`. The `GRangesList` contains the Loci of interest. Finally, the individual `GRanges` contains ChIP-seq like scores for every n base pairs (with $n = \text{stepSize}$, see below).

This object may be difficult to navigate if many different parameters, or Loci are used. In order to facilitate navigation, we included a search function. **See function: `searchSites`** This function can also be used to navigate `AllSitesAboveThreshold` slot after occupancy scores have been computed.

Step 6 - Model Accuracy

In order to plot the model accuracy (predicted model against real ChIP-seq data).

```
AccuracyEstimate <- profileAccuracyEstimate(LocusProfile = eveLocusChip,
  predictedProfile = chipProfile, occupancyProfileParameters = OPP)
AccuracyEstimate
```

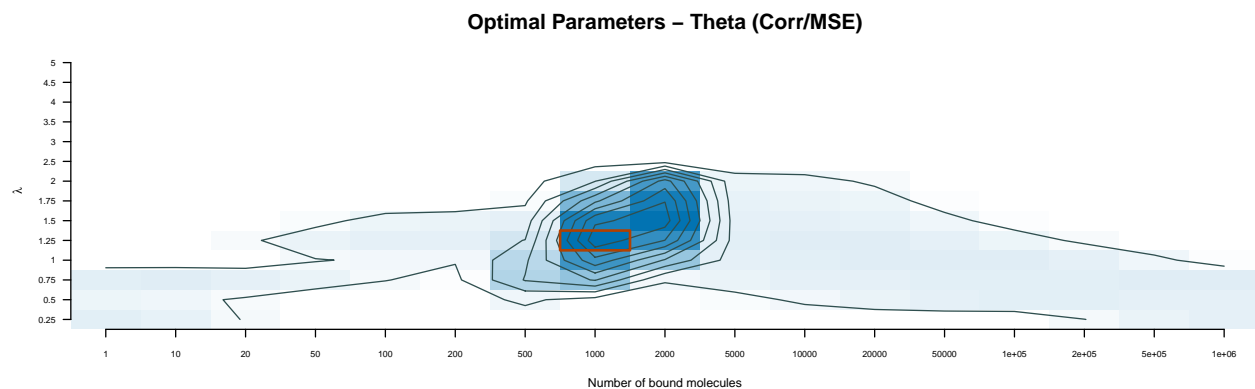
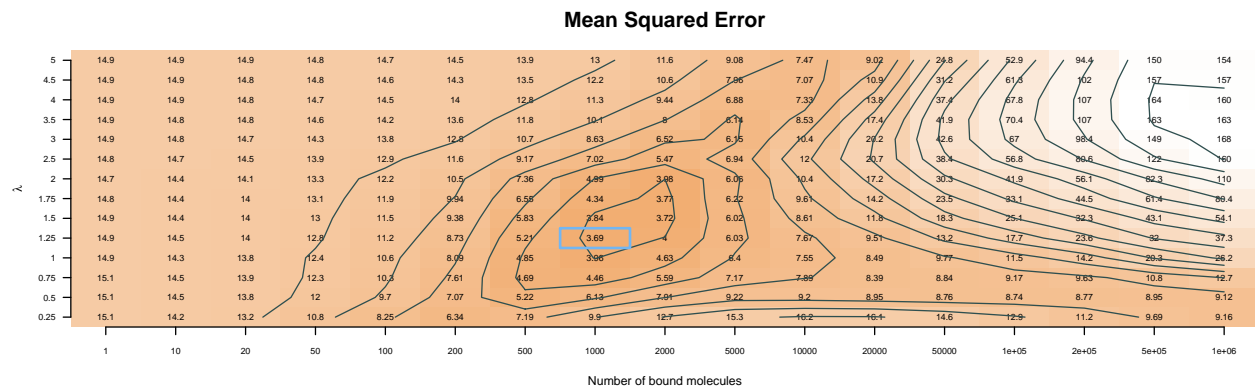
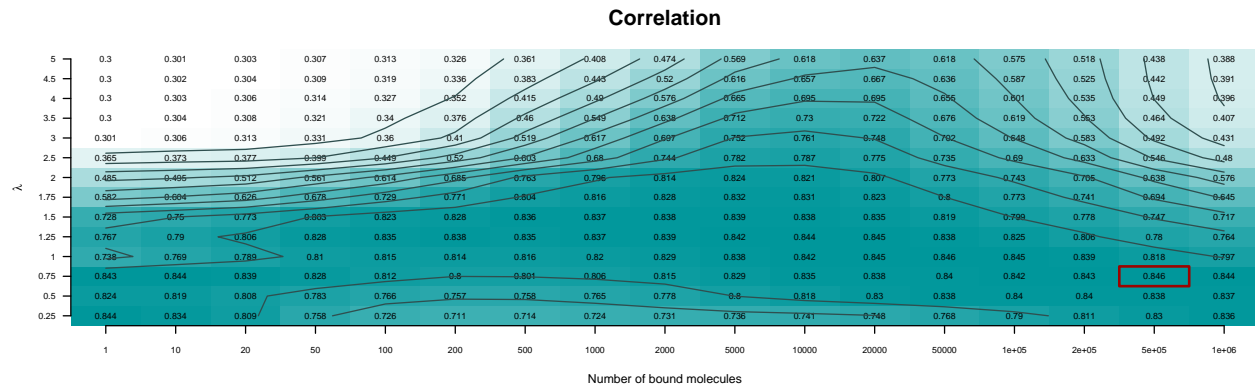
```
## $`lambda` = 1.5 & boundMolecules = 1000`
## $`lambda` = 1.5 & boundMolecules = 1000`$eve
```

```
##          Corr          MSE    meanCorr    meanMSE    meanTheta
## 0.836843223 0.003835814 0.836843223 3.835814393 0.218165724
```

Step 7 - Plotting

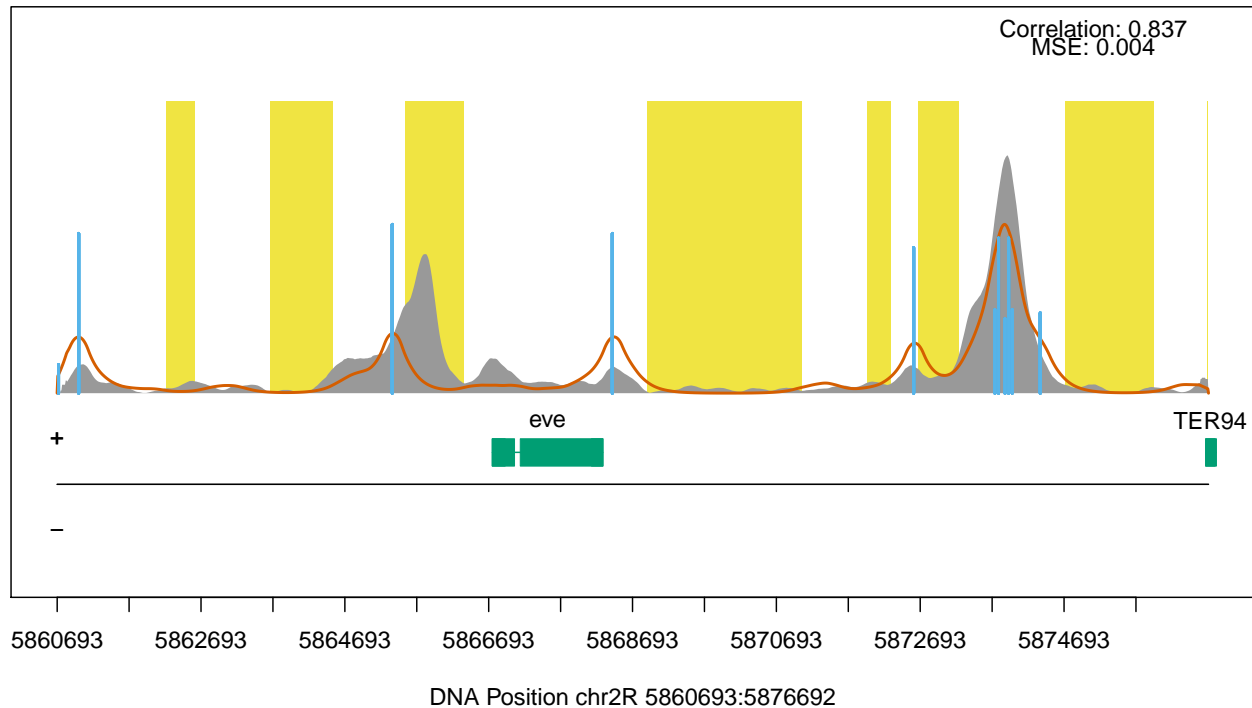
Finally, once all has been computed, it is possible to plot the results.

```
# Plotting Optimal heat maps
plotOptimalHeatMaps(optimalParam, parameter="all")
```



```
# Plotting occupancy Profile
plotOccupancyProfile(predictedProfile=chipProfile[[1]][[1]],
  setSequence=eveLocus,
  profileAccuracy = AccuracyEstimate[[1]][[1]],
  chipProfile = eveLocusChip[[1]],
  occupancy = AllSitesAboveThreshold(Occupancy)[[1]][[1]],
```

```
DNAAccessibility = Access,
occupancyProfileParameters = OPP,
geneRef = geneRef)
```



Work Flow - Full Guide

This section will describe ChIPAnalyser's work flow. However in this section we will describe in detail data objects, parameters, and functions. Please refer to this section if in doubt. If the doubt persists, don't hesitate to send an email to the maintainer.

Data objects - Genomic Profile Parameters

The very first aspect to consider when using ChIPAnalyser is data input. Many (if not all functions) require specific data inputs and parameters in order to carry out the computation. To facilitate, the storage of these parameters, we created a **genomicProfileParameters** object (S4 class). This is the very first step before any other work. All other functions rely on this **genomicProfileParameters** object in one form or another. The output of most functions will be a **genomicProfileParameters** object. Thus the output of one functions should be used as an input for the next functions in the pipeline. All functions are described bellow in section **Work Flow - Analysis**.

This object comes in the following form:

```
genomicProfileParameters(PWM, PFM, ScalingFactorPWM, PFMFormat, pseudocount,
  BPFrequency, naturalLog, noOfSites,
  minPWMScore, maxPWMScore, PWMThreshold,
  AllSitesAboveThreshold, DNASequenceLength,
  averageExpPWMScore, strandRule, whichstrand, NoAccess)
```

To build a **genomicProfileParameters** object :


```
# Assign Value wanted for each parameter
GPP <- genomicProfileParameters(PWM, PFM,ScalingFactorPWM, PFMFormat,
  pseudocount, BPFfrequency, naturalLog, noOfSites,
  PWMThreshold, DNASequencLength,
  strandRule, whichstrand)
```

As one can see, `genomicProfileParameters` contains many arguments. However many of these arguments already have default values assigned to them. Some of the arguments should not be set by user. These values are computed internally and will automatically updated (`minPWMScore`, `maxPWMScore`, `AllSitesAboveThreshold`, `NoAccess`). In this situation, most arguments are not required to build a `genomicProfileParameters` object and a minimal build can be described as:

```
# return empty genomicProfileParameters object
GPP <- genomicProfileParameters()
# return minimal working object
GPP <- genomicProfileParameters(PFM=PFM,PFMFormat="raw")
# Suggested Minimal Build
GPP <- genomicProfileParameters(PFM=PFM,PFMFormat="raw",
  BPFfrequency=DNASequencSet)
```

Although many parameters have assigned default values, it is recommended to use custom parameters to better fit the needs of the analysis. The method described above will build a new `genomicProfileParameters` object with the values that were assigned to each argument. Only three slots are required in order to build a `genomicProfileParameters` object (see below - **The compulsory ones**). Most other slots are optional. If after building `genomicProfileParameters`, you wish to modify the value of only *one* slot and keep the values that you had previously assigned, it is possible to modify each slot individually by using the slot *access/setter* methods. Each slot and it's *access/setter* method is described below.

Position Matrices - The compulsory ones

- **PWM** , a Position Weight Matrix. If a Position Weight Matrix is readily available it is possible to directly use this Matrix. This PWM should contain four rows (one for each base pair; ACTG in order). The number c olumns will depend on the length of the preferred binding motif of a given Transcription Factor. This argument is only necessary IF and ONLY IF, no PFM (Position Frequency Matrix) is available. Choosing between PWM or PFM comes down to personal choice as long a PWM is available for further computation (see PFM). If a PFM is available (see below), the Position Weight Matrix will be directly computed from the Position Frequency Matrix. Although it is possible to assign a new PWM to the `genomicProfileParameters` object without creating a new object, we suggest that if you were to decided to use another Position Weight Matrix to create a new `genomicProfileParameters`.

```
#Accessing PositionWeightMatrix slot
PositionWeightMatrix(GPP)
```

```
##      [,1]      [,2]      [,3]      [,4]      [,5]      [,6]      [,7]
## A  0.1267378 -0.8713677 -3.953162  1.983869  1.983869 -5.052697 -9.445015
## C  0.2913871  0.6224195 -4.801159 -9.445015 -9.445015 -9.445015  1.998447
## G  0.3703684 -2.3054635 -9.445015 -4.587034 -4.587034 -3.422647 -9.445015
## T -1.3522577  0.7753635  1.962784 -9.445015 -9.445015  1.954263 -9.445015
##      [,8]
## A -4.235561
## C  1.831691
## G -3.830305
## T -1.657112
```

```
# Setting PositionWeightMatrix slot
PositionWeightMatrix(GPP) <- newPWM
### This is not the advised method
### newPWM is a matrix following the format described above
```

- PFM, a Position Frequency Matrix. The Position Frequency Matrix argument may come in multiple forms: in the form of a Matrix containing four rows (one for each base pair ACTG) and columns depending of the length of the binding motif or in the form of a path to file linking to a PFM. Position Frequency Matrices come in various configurations. The most common ones (all supported by ChIPAnalyser) are RAW (similar to the simple matrix described previously), Transfac and JASPAR. Finally, if the binding sequences are available, the PFM will be generated from sequence information. We suggest to use a path/to/file linking towards the PFM file. Most PFM will come in one of the formats described above and ChIPAnalyser will parse these files in a usable format. However, PLEASE NOTE THAT THE FORMAT SHOULD BE SPECIFIED. See PFMFormat below.

If a PWM is readily available, PFM is not necessary. However, keep in mind that at least one is necessary. Although it is possible to assign a new PFM to the `genomicProfileParameters` object without creating a new object, we suggest that if you were to decided to use another Position Frequency Matrix to create a new `genomicProfileParameters`.

```
# Accessing PositionFrequencyMatrix slot
PositionFrequencyMatrix(GPP)
```

```
##      [,1] [,2] [,3] [,4] [,5] [,6] [,7] [,8]
## A   190   95   11  689  689    5    0    9
## C   213  268    6    0    0    0  696  620
## G   225   35    0    7    7   16    0   12
## T    68  298  679    0    0  675    0   55
```

```
# Setting PositionFrequencyMatrix slot
PositionFrequencyMatrix(GPP) <- newPFM
```

In this situation, `newPFM` is either a path to file or a PFM matrix. The `PFMFormat` will be the one assigned to the `genomicProfileParameters` object.

At least one of **PWM** or **PFM** is required to create a `genomicProfileParameters` storage object. If a PFM is provided then the PWM will be automatically computed and updated.

- `PFMFormat`, a file format for `PositionFrequencyMatrix` file. When Loading a PFM from a file (as described above), one should included the format of the file that they are using. `PFMFormat` may be one of the following: “raw”, “transfac”, “JASPAR” or “sequences”.

```
PFMFormat(GPP)
```

```
PFMFormat(GPP)<-"raw"
```

Default is set at “raw”.

All other arguments are optional however we strongly recommend to tailor the values assigned to `genomicProfileParameters` to your needs. The following sections will describe these optional parameters.

Genomic Parameters - The optional ones

- `ScalingFactorPWM`, a scaling factor for TF specificity. Although this parameter is optional (Default value is set at 1), the *scaling factor* (or *lambda* as described in the equations above) is crucial for many functions (described below). `ScalingFactorPWM`, must be a positive numeric value or a vector containing positive numeric values. The optimal value for `ScalingFactorPWM` may be inferred by using

`computeOptimal`. Different values for `ScalingFactorPWM` will influence the goodness of fit of the model. For more information, see `computeOptimal` and `profileAccuracyEstimate`.

```
ScalingFactorPWM(GPP)
```

```
ScalingFactorPWM(GPP) <- 0.5
```

```
ScalingFactorPWM(GPP) <- c(0.5, 1, 1.5, 2)
```

- **PWMPseudocount**, a probability modifier. When computing a PWM from a PFM, it is possible that certain base pairs are completely absent from the Position Frequency Matrix. This absence will lead to odd results as part of this transformation requires a logarithmic transformation (at Position probability matrix step - a Matrix that describes the simple probability of a base pair being in that position of a binding motif given the PFM). *zeroes* will give minus infinities. In order to overcome this problem, a **PWMPseudocount** is introduced in the Position Probability Matrix. a **PWMPseudocount** of 1 (Default Value is 1) will then become a 0 after logarithmic transformation thus removing any mathematical discomforts.

```
PWMPseudocount(GPP)
```

```
PWMPseudocount(GPP) <- 1
```

- **BPFrequency**, the frequency at which each base pair will occur in a given organism. Probabilistically speaking, all base pairs have an equal chance of occurring in the genome (Default value for this slot is set at 0.25 per base pair). However, biologically speaking this is not the case. **BPFrequency** may be supplied in various forms. If base pair frequency is known, it may be supplied as a vector containing the probability of occurrence of each base pair. If however, this frequency is unknown, **genomicProfileParameters** will compute **BPFrequency** from a **BSgenome** or a **DNAStringSet**. Bare in mind that **BPFrequency** is used to generate a *PWM* from a *PFM*, thus if one were to change the **BPFrequency** after creating a **genomicProfileParameters** with an already computed *PWM*, this would not influence the value of the *PWM*. It would be necessary to rebuild a new **genomicProfileParameters** object.

```
BPFrequency(GPP)
```

```
BPFrequency(GPP) <- c(0.2900342, 0.2101426, 0.2099192, 0.2899039)
```

```
BPFrequency(GPP) <- DNASequenceSet
```

- **naturalLog**, a logical value. As described previously (see **pseudocount**), the transformation from PFM to PWM requires a logarithmic transformation. The user may choose which logarithmic transformation, they would rather apply (Default is **TRUE**). If **naturalLog** = **TRUE**, then the natural logarithm will be used for transformation. If **naturalLog** = **FALSE**, then *log2* will be used instead. Keep in mind that, the goal is to avoid any funky business during PFM to PWM transformation (e.g. Minus infinities or division by zero).

```
naturalLog(GPP)
```

```
naturalLog(GPP) <- FALSE
```

- **noOfSites**, the number of sites used to compute the PWM from the PFM. In the event that a PFM contains a large amount of sites (as it sometimes is the case with Transfac PFM), it is possible to restrict this number of sites. The default value is 0. When **noOfSites** = 0, the whole PFM is used to compute the PWM.

```
noOfSites(GPP)
```

```
noOfSites(GPP) <- 8
```

- **PWMThreshold**, a numeric threshold against which PWM Scores are selected (Default is 0.7). Although it is possible to compute every single motif present in a stretch of DNA (if this is of interest, set **PWMThreshold** to 0), in most cases, only the sites with a high PWM Score will be of interest. The **PWMThreshold**, a numeric value between 0 and 1, will select regions above that given threshold. For the default threshold of 0.7, only the top 30% of PWM Scores will be selected.

```
PWMThreshold(GPP)
```

```
PWMThreshold(GPP) <- 0.7
```

- **strandRule**, indicates how the genome should be scored with the PWM (Default is "max"). As DNA is double stranded, it is necessary to specify how a strand of DNA should be scored. If **strandRule** = "max", both strands will be scored and the highest score between each strand will be selected. If **strandRule** = "sum", both strands will be scored and their respective score will be summed. If **strandRule** = "mean", both strands will be scored and the average score between both strands will be selected as PWM Score. Only three possibilities: "max", "sum" and "mean"

```
strandRule(GPP)
```

```
strandRule(GPP) <- "mean"
```

- **whichstrand**, indicates which strand will be used to score the genome with the PWM (Default is both strand and is indicated by "+-"). Three options exist: plus strand ("+"), minus strand ("-") or both ("+-" or "-+").

```
whichstrand(GPP)
```

```
whichstrand(GPP) <- "+"
```

Genomic Parameters - The Updated ones

Some of the slots **genomicProfileParameters** should not be changed by user. We strongly advise against changing these slots. Certain Parameters are updated after a certain computation has been carried out. For example, **maxPWMScore** and **minPWMScore** are computed during the **computeGenomeWidePWMScore** function (see below) and represent both the highest and the lowest score of the given DNA sequence. These slots will be updated in the **genomicProfileParameters** object as one makes its way through the ChIPAnalyser work flow. Essentially, they are place holders for information required further down the work flow. Only slots that are of interest for the user are available for visualisation. If these slots have not been updated, the function will not return any value.

- **maxPWMScore**, a numeric value describing the highest PWM Score on a given DNA sequence and the value assigned to **lambda**. It is still possible to access this slot using:

```
maxPWMScore(Occupancy)
```

```
## [1] 12.86065
```

- **minPWMScore**, a numeric value describing the lowest PWM Score on a given DNA sequence and the value assigned to **lambda**. It is possible to access this slot using:

```
minPWMScore(Occupancy)
```

```
## [1] -48.82628
```

- **averageExpPWMScore** a numeric value representing the exponential of the average PWM Score. This score depends on the values assigned to **lambda**. It is possible to access this slot using:

```
averageExpPWMScore(Occupancy)
```

```
## [1] 1.015637
```

- **DNASequenceLength**, a numeric value describing the length of the DNA sequence used. Although theoretically one could provide this information, DNA length is automatically computed and the slot updated during **computeGenomeWidePWMScore** function. The length of this sequence is the length of the sequence used to compute the scores previously mentioned (**maxPWMScore**, **minPWMScore** and **averageExpPWMScore**). This means that if DNA accessibility data is provided, the length of the sequence will only be the length of the accessible DNA.

```
DNASequenceLength(Occupancy)
```

```
## [1] 3112514
```

- **NoAccess**, indicates if certain Loci of interest (see **setSequence** below) **do not** contain any accessible DNA. It is possible that certain of the loci you have chosen do not contain any accessible DNA (no overlap with DNA accessibility data provided). If this is the case, you will be notified during the computation and the loci will be stored in the **NoAccess** slot.

```
NoAccess(Occupancy)
```

```
## [1] "-"
```

- **AllSitesAboveThreshold**, stores all sites above threshold with the associated PWM Score and Occupancy. This slot may contain a variety of objects however they all represent the same thing: it will always contain at its core a **GRanges** object (slot class defined as "GRlist" - can be one of the following **GRangesList** or **list**). This **GRanges** includes sites above threshold (start, end and strand), **PWMScores** for those sites and possibly **Occupancy** (depending on what has already been computed). **GRanges** are encapsulated in a **GRangesList** as each **GRanges** represent a specific Loci. This **GRangesList** may also be encapsulated in a list. This list will represent a combination of **lambda** and number of bound Molecules (see **boundMolecules**). For more information on this list see **computeOccupancy**. It is possible to access this slot by using:

```
AllSitesAboveThreshold(Occupancy)
```

```
## $`lambda` = 1.5 & boundMolecules = 1000`
## GRangesList object of length 1:
## $eve
## GRanges object with 412 ranges and 2 metadata columns:
##      seqnames      ranges strand |      PWMScore
##      <Rle>        <IRanges> <Rle> |      <numeric>
## eve chr2R [5860705, 5860712]   + | -1.51655573585429
## eve chr2R [5860709, 5860716]   + | -5.33217184502491
## eve chr2R [5860715, 5860722]   + |  9.13992557549757
## eve chr2R [5860728, 5860735]   + |  5.05434682102833
## eve chr2R [5860758, 5860765]   + | -5.15370980167748
## ...      ...                ...   ...
## eve chr2R [5876629, 5876636]   + |  5.60817413411963
## eve chr2R [5876635, 5876642]   + |  0.202790199774102
## eve chr2R [5876641, 5876648]   - | -4.47385601266488
## eve chr2R [5876666, 5876673]   + |  2.21133362723558
## eve chr2R [5876684, 5876691]   + | -2.28895797651261
##      Occupancy
##      <numeric>
## eve 0.0138683203024566
## eve 0.0138160293072631
## eve 0.0783704718441574
## eve 0.0183246335750422
## eve 0.0138165926852252
```

```
## ...
## eve 0.0203268664876583
## eve 0.0139901018008407
## eve 0.0138194725807681
## eve 0.014492381648295
## eve 0.0138454813805688
##
## -----
## seqinfo: 1 sequence from an unspecified genome; no seqlengths
```

Or

```
searchSites(Occupancy)
```

```
## $`lambda` = 1.5 & boundMolecules = 1000`
## GRangesList object of length 1:
## $eve
## GRanges object with 412 ranges and 2 metadata columns:
##      seqnames      ranges strand |      PWMScore
##      <Rle>         <IRanges> <Rle> |      <numeric>
## eve chr2R [5860705, 5860712]   + | -1.51655573585429
## eve chr2R [5860709, 5860716]   + | -5.33217184502491
## eve chr2R [5860715, 5860722]   + |  9.13992557549757
## eve chr2R [5860728, 5860735]   + |  5.05434682102833
## eve chr2R [5860758, 5860765]   + | -5.15370980167748
## ...      ...      ...      ... |      ...
## eve chr2R [5876629, 5876636]   + |  5.60817413411963
## eve chr2R [5876635, 5876642]   + |  0.202790199774102
## eve chr2R [5876641, 5876648]   - | -4.47385601266488
## eve chr2R [5876666, 5876673]   + |  2.21133362723558
## eve chr2R [5876684, 5876691]   + | -2.28895797651261
##      Occupancy
##      <numeric>
## eve 0.0138683203024566
## eve 0.0138160293072631
## eve 0.0783704718441574
## eve 0.0183246335750422
## eve 0.0138165926852252
## ...      ...
## eve 0.0203268664876583
## eve 0.0139901018008407
## eve 0.0138194725807681
## eve 0.014492381648295
## eve 0.0138454813805688
##
## -----
## seqinfo: 1 sequence from an unspecified genome; no seqlengths
```

The size of the `AllSitesAboveThreshold` slot will increase drastically as the number of values assigned to `ScalingFactorPWM` (or `lambda`) and `boundMolecules` increases. In order to navigate and search this slot with ease, it is possible to use the `searchSites` function (See below: `searchSites`).

Data Objects - Occupancy Profile Parameters

`genomicProfileParameters` represents a good chunk of the parameters needed to go through the entire ChIPAnalyser work flow. However, there are more to come! A second parameter storing object was created to handle non-compulsory parameters. This lightens `genomicProfileParameters` by handling part of the parameters. This second S4 object is called `occupancyProfileParameters`. The interesting aspect of this object is that none of the slots are compulsory. This means that if not provided, a new `occupancyProfileParameters` object will be created internally. All default values will be used for further computation. As stated previously, we strongly advise using custom parameters in order to increase goodness of fit of model. It is especially the case here, as slots such as `maxSignal` are directly extracted from biological data (ChIP-seq data - see `computeChIPProfile` and `profileAccuracyEstimate` for more information).

```
OPP <- occupancyProfileParameters(ploidy = 2 ,boundMolecules = 1000 ,
  backgroundSignal = 0 ,maxSignal = 1, chipMean = 150 , chipSd = 150 ,
  chipSmooth = 250 , stepSize = 10 ,
  removeBackground = 0 , thetaThreshold = 0.1)
```

As it is the case with `genomicProfileParameters`, it is also possible to *access/set* each slot individually after having created an `occupancyProfileParameters` object. Each slot is described as the following:

- `ploidy`, the ploidy level of the organism of interest (Default is set at 2). This only considers simple polyploidy (or haploidy). The model does not (yet) consider hybrids such as wheat.

```
ploidy(OPP)
ploidy(OPP) <- 2
```

- `boundMolecules`, a positive integer (or vector of positive integers) describing the number of bound molecules (Transcription factors) to DNA (Default value is set at 2000). In this model, occupancy is reliant on the number of bound molecules. The number of molecules will influence the goodness of fit of the model. It is possible to infer the number of bound Molecules by using the `computeOptimal` function. For more information, see `computeOptimal` and `profileAccuracyEstimate`.

```
boundMolecules(OPP)
boundMolecules(OPP) <- 5000
```

- `backgroundSignal`, a numeric value representing the background Signal in real ChIP-seq data (Default is set at 0). It is strongly advised to set this parameter to the background Signal of the ChIP-seq data you will be using.

```
backgroundSignal(OPP)
backgroundSignal(OPP) <- 0.02550997
```

- `maxSignal`, a numeric value representing the maximum signal in real ChIP-seq data (Default is set at 1). It is strongly advised to set this parameter to the maximum Signal of the ChIP-seq data you will be using.

```
maxSignal(OPP)
maxSignal(OPP) <- 1.86
```

- `chipMean`, a numeric value representing the average peak width in base pairs in real ChIP-seq data (Default is set at 150). It is strongly advised to set this parameter to the average peak width of the ChIP-seq data you will be using.

```
chipMean(OPP)
chipMean(OPP) <- 150
```


- **chipSd**, a numeric value representing the standard deviation of peak width in real ChIP-seq data (Default is set at 150). It is strongly advised to set this parameter to the SD peak width of the ChIP-seq data you will be using.

```
chipSd(OPP)
```

```
chipSd(OPP) <- 150
```

- **chipSmooth**, a numeric value representing the size of the window used for smoothing the profile (Default is set at 250). The goal of ChIPAnalyser is to produce ChIP-seq like profile from predicted high occupancy sites. In order to mimic these ChIP-seq profile, a smoothing algorithm is used to smooth occupancy profiles. This algorithm uses ChIP-seq parameters such as **chipMean**, **chipSd**, **maxSignal**, **backgroundSignal** and **chipSmooth**.

```
chipSmooth(OPP)
```

```
chipSmooth(OPP) <- 250
```

- **stepSize**, a numeric value describing the bin size (in base pairs) used for computing ChIP-seq like profiles (Default is set at 10). In the case of long sequences, it not always necessary to include ChIP-like occupancy at every base pair (mainly for speed and memory usage). **stepSize** will determine the size of the bins used to split your sequence of interest. As an example, if your sequence is 16 000 bp long with a **stepSize** of 10, the resulting profile will be composed of 1600 occupancy points.

```
stepSize(OPP)
```

```
stepSize(OPP) <- 10
```

- **removeBackground**, a numeric value describing a threshold at which Occupancy signals must be removed (Default is set at 0).

```
removeBackground(OPP)
```

```
removeBackground(OPP) <- 0
```

- **thetaThreshold**, a numeric value describing the threshold used to calculate our in house *theta* value (Default is set at 0.1). *Theta* is a metric used to demonstrate which parameters are optimal by maximising the correlation and minimising the Mean Squared Error (MSE) between the predicted profile and actual ChIP-seq profiles. The higher the value of *theta*, the better the ratio between correlation and MSE. Values below this threshold are discarded (replaced by Threshold) as they represent extremely poor accuracy with actual ChIP-seq data.

```
thetaThreshold(OPP)
```

```
thetaThreshold(OPP) <- 0.1
```

Work Flow - Analysis

Once a **genomicProfileParameter** object has been established, the rest of the analysis becomes fairly straight forward. Unless, you already have prior knowledge on the number of bound molecules (**boundMolecules**) and the PWM scaling factor (**ScalingFactorPWM** or referred to as *lambda*), we advise you to first infer the optimal set of parameters as described in **computeOptimal**. However, as this function is essentially a combination of all other functions in the package (with a little bit more magic to it), we will overview a simple analysis work flow first and finish with **computeOptimal** function and its associated plotting function **plotOptimalHeatMaps**.

Genome Wide Scoring

In order to score the entire genome (or the accessible genome), it is possible to use the `computeGenomeWidePWMScore` function. The output of this function will be influenced by the value assigned to `lambda`. If more than one value was assigned to the scaling factor, parameters dependant on `lambda` will be updated accordingly (computed for each value of `lambda`). The arguments of the function are the following :

```
computeGenomeWidePWMScore(DNASequenceSet, genomicProfileParameters,  
  DNAAccessibility = NULL, GenomeWide = TRUE, verbose = TRUE)
```

Input Data - Genome Wide scoring

As input, `computeGenomeWidePWMScore` requires to obligatory arguments: `DNASequenceSet` and `genomicProfileParameters`. `DNASequenceSet` comes in the form of the following:

`DNASequenceSet`

```
## A DNAStringSet instance of length 15  
##      width seq                                     names  
## [1] 23011544 CGACAATGCACGACAGAGG...ATGAACCCCCCTTTCAAA chr2L  
## [2] 21146708 GACCCGCTAGGAGATGTTG...TTTGCATTCTAGGAATTC chr2R  
## [3] 24543557 TAGGGAGAAATATGATCGC...AACCAAGTTAATGTTTCGG chr3L  
## [4] 27905053 GAATTCTCTCTTGTGTAG...TTCGCATTCTAGGAATTC chr3R  
## [5] 1351857 GAATTGCGCTCCGCTTACC...CGATTTGAGATATATGAA chr4  
## ...  
## [11] 2555491 AACGAGGCCCATTTTCATAC...ATGCCATTGCTAGAAGT chr3LHet  
## [12] 2517507 CCCTGTTTGCATCAGCGTT...TAAAAACAATTTGCTCCC chr3RHet  
## [13] 204112 TAGATAGATAGATAGATAG...ATCGGAGTTAATGTTTGC chrXHet  
## [14] 347038 AGGGTCACGTAATGCTGAT...TTGTTTCCCCGGGATTG chrYHet  
## [15] 29004656 ATTGAAATGGATTGCATT...CAAGACCTTTCAAGACAA chrUextra
```

`DNASequenceSet` may also come in the form of a `BSgenome` object. However, we advise to use a `DNAStringSet` for a question of ease and speed. If you are unfamiliar with `BSgenome` and `DNAStringSet`, the following example demonstrates how to use these objects in this context.

```
#Extracting DNAStringSet from BSgenome
```

```
DNASequenceSet <- getSeq(BSgenome.Dmelanogaster.UCSC.dm3)
```

As a reminder a `genomicProfileParameters` are presented in the following format:

`GPP`

```
## Object Class:genomicProfileParameters
```

```
##
```

```
##
```

```
## PWM:
```

```
##      [,1]      [,2]      [,3]      [,4]      [,5]      [,6]      [,7]  
## A  0.1267378 -0.8713677 -3.953162  1.983869  1.983869 -5.052697 -9.445015  
## C  0.2913871  0.6224195 -4.801159 -9.445015 -9.445015 -9.445015  1.998447  
## G  0.3703684 -2.3054635 -9.445015 -4.587034 -4.587034 -3.422647 -9.445015  
## T -1.3522577  0.7753635  1.962784 -9.445015 -9.445015  1.954263 -9.445015  
##      [,8]  
## A -4.235561  
## C  1.831691  
## G -3.830305
```

```

## T -1.657112
##
## PFM:
##   [,1] [,2] [,3] [,4] [,5] [,6] [,7] [,8]
## A  190   95   11  689  689    5    0    9
## C  213  268    6    0    0    0  696  620
## G  225   35    0    7    7   16    0   12
## T   68  298  679    0    0  675    0   55
##
## PFMFormat: raw
##
## PWM Scores at Sites higher than Threshold:
## GRangesList object of length 0:
## <0 elements>
##
## -----
## seqinfo: no sequences
##
## No Accessible DNA at Loci:
##
## Genomic Profile Parameters:
## Lambda: 1
## BP Frequency:    0.25    0.25    0.25    0.25
## Pseudocount: 1
## Natural log: FALSE
## Number Of Sites: 0
## maxPWMScore:
## minPWMScore:
## PWMThreshold: 0.7
## Average Exponential PWM Score:
## DNA Sequence Length:
## Strand Rule: max
## Strand: +-

```

DNAAccessibility is an optional argument in `computeGenomeWidePWMScore`. If present, then the genome will be scored only on the accessible DNA. DNAAccessibility comes as a `GRanges` containing accessible DNA sites.

```

# DNA accessibility
Access

```

```

## GRanges object with 4703 ranges and 0 metadata columns:
##           seqnames           ranges strand
##           <Rle>             <IRanges> <Rle>
##      [1]   chr2R [ 7339296,  7342564]   *
##      [2]   chr2R [ 9436993,  9437589]   *
##      [3]   chr2R [15728083, 15728687]   *
##      [4]   chr2R [ 4980200,  4980845]   *
##      [5]   chr2R [ 6028863,  6029419]   *
##      ...      ...                ...   ...

```

```
## [4699] chr2R [21120053, 21120400] *
## [4700] chr2R [21140572, 21140980] *
## [4701] chr2R [21143160, 21143517] *
## [4702] chr2R [21144932, 21145281] *
## [4703] chr2R [21145564, 21146702] *
## -----
## seqinfo: 6 sequences from an unspecified genome; no seqlengths
```

Finally, `verbose` will determine if progress messages should be printed in the console.

computeGenomeWidePWMScore

As an example of `computeGenomeWidePWMScore` usage:

```
# With DNAAccessibility

GenomeWide <- computeGenomeWidePWMScore(DNASequenceSet = DNASequenceSet,
  genomicProfileParameters = GPP, DNAAccessibility = Access)

GenomeWide

# Without DNA accessibility

GenomeWide <- computeGenomeWidePWMScore(DNASequenceSet = DNASequenceSet,
  genomicProfileParameters = GPP)

GenomeWide
```

Scoring sites above threshold

Once genome wide metrics have been computed, the next step in the analysis is to extract sites above threshold (Sites with strong binding sites according to PWM Scores). The `computePWMScore` function will score the genome and extract sites above a local threshold (dependant on `PWMThreshold`, `maxPWMScore` and `minPWMScore`). The arguments of this functions are the following:

```
computePWMScore(DNASequenceSet, genomicProfileParameter,
  setSequence = NULL, DNAAccessibility = NULL, verbose = TRUE)
```

Input Data - Sites Above threshold

Only two arguments are absolutely required: `DNASequenceSet` and `genomicProfileParameters`. However, `setSequence` represents the Loci of interest. If `setSequence = NULL`, then sites above threshold will be computed and extracted on a genome wide scale (or accessible genome if DNA Accessibility is provided). `DNASequenceSet` and `DNAAccessibility` are in the same format as previously described (`verbose` plays the same role as previously described). `setSequence` is a `GRanges` representing the loci of interest (may contain more than one loci/range) and comes in the following format:

```
eveLocus

## GRanges object with 1 range and 0 metadata columns:
##      seqnames      ranges strand
##      <Rle>         <IRanges> <Rle>
## eve chr2R [5860693, 5876692] *
## -----
## seqinfo: 1 sequence from an unspecified genome; no seqlengths
```

An important aspect to mention, is that it is imperative you name your loci of interest (not to be confused with `seqnames`). If you are unfamiliar with `GRanges`, the following examples demonstrates naming in the context of `ChIPAnalyser`. We recommend getting acquainted with `GenomicRanges` as many aspect of `ChIPAnalyser` require the use of `GRanges`.

```
# Sequence names of Loci
seqnames(eveLocus)

## factor-Rle of length 1 with 1 run
##   Lengths:      1
##   Values : chr2R
## Levels(1): chr2R

# Names of Loci
names(eveLocus)

## [1] "eve"

# Naming Loci in GRanges
names(eveLocus) <- "eve"
```

computePWMScore

To compute PWM Scores at sites above threshold:

```
# With DNA Accessibility

PWMScores <- computePWMScore(DNASequenceSet = DNASequenceSet,
                             genomicProfileParameters = GenomeWide,
                             setSequence = eveLocus, DNAAccessibility = Access)
PWMScores

# Without DNA Accessibility

PWMScores <- computePWMScore(DNASequenceSet = DNASequenceSet,
                             genomicProfileParameters = GenomeWide,
                             setSequence = eveLocus)
PWMScores
```

As you can see, the `genomicProfileParameters` argument is the `genomicProfileParameters` object computed in the previous example. `ChIPAnalyser` works in a sequential manner: resulting object from one functions are often parsed as arguments to other functions. Finally, if your sequence of interest does not contain any accessible DNA, you will be notified during the computation and it is possible to extract inaccessible loci by using `NoAccess(PWMScores)` (See `NoAccess` slot in `genomicProfileParameters`).

Occupancy

Occupancy scores are computed using the formula described in **Methods**. It is worth mentioning that Occupancy scores are dependant on values assigned to `ScalingFactorPWM` and `boundMolecules`. If more than one value were to be assigned to these parameters, the resulting output will be a combination of both. For more information see the `computeOccupancy` example as we will demonstrate multiple value computation (Single Value for `lambda` and `boundMolecules` will return an object identical in structure as with multiple values). The arguments for `computeOccupancy` are the following:

```
computeOccupancy(AllSitesPWMScore, occupancyProfileParameters = NULL,
  norm = TRUE, verbose = TRUE)
```

Input Data - Occupancy

`computeOccupancy` requires a `genomicProfileParameters` object result of the previous function (`computePWMScore`). If you are unsure, if your `genomicProfileParameter` contains the right information, it is possible to check by using:

```
AllSitesAboveThreshold(PWMScores)
```

If your `GRanges` does not contain `PWMScore` as a metadata column, you are either using the wrong object or you have not yet computed PWM Scores.

`occupancyProfileParameters` is an `occupancyProfileParameters` object. If not provided, a new one will be generated internally. As previously mentioned, we strongly recommend to set those parameters to improve the model's goodness of fit. As a reminder, a `occupancyProfileParameters` object (previously created - see section **Data object - Occupancy profile Parameters**) should print on the screen as follows:

```
OPP
```

```
## Object Class:occupancyProfileParameters
##
## Ploidy: 2
## boundMolecules: 1000
## backgroundSignal: 0.02550997
## maxSignal: 1.847
## chipMean: 200
## chipSd: 200
## chipSmooth: 250
## Step Size: 10
## Theta Threshold: 0.1
```

Finally, if `norm = TRUE`, the occupancy profiles will be normalised and `verbose = TRUE` progress messages will be printed to the console.

computeOccupancy

To compute Occupancy scores with `computeOccupancy`:

```
Occupancy <- computeOccupancy(AllSitesPWMScore = PWMScores,
  occupancyProfileParameters = OPP)
Occupancy
```

As it is the case in the previous functions, `AllSitesPWMScore` should be the result of the previous function (`computePWMScore`). `computeOccupancy` will return a `genomicProfileParameters` object with an updated `AllSitesAboveThreshold` slot. This slot should now contain a list of `GRangesLists` containing `GRanges` (one for each Loci of interest) with two metadata columns (`PWMScore` and `Occupancy`). Each element in the list is named with the specific combination of *lambda* and *boundMolecules* used to compute this set of occupancies. Finally, if your sequence of interest does not contain any accessible DNA, you will be notified during the computation and it is possible to extract inaccessible loci by using `NoAccess(PWMScores)` (See `NoAccess` slot in `genomicProfileParameters`).

ChIP-seq like profiles

The ultimate goal of ChIPAnalyser is to produce *ChIP-seq like* profile from occupancy data (from sites that display a high TF occupancy). `computeChipProfile` creates *ChIP-seq like* profiles from occupancy data by smoothing occupancy *profiles* and mimicking real ChIP-seq data. The arguments of `computeChipProfile` are the following:

```
computeChipProfile( setSequence ,
  occupancy, occupancyProfileParameters = NULL, norm = TRUE,
  method = c("moving_kernel","truncated_kernel","exact"),
  peakSignificantThreshold= NULL,
  verbose = TRUE)
```

Input data - ChIP-seq profiles

The `computeChipProfile` function requires two compulsory arguments `setSequence` and `occupancy`. `setSequence` is a `GRanges` describing the loci of interest (this is the same `GRanges` used in `computePWMScore`). `occupancy` is a `genomicProfileParameters` object result of `computeOccupancy` function. To make sure this is the right `genomicProfileParameters`, you may use `AllSitesAboveThreshold()` (See `AllSitesAboveThreshold` slot description above). `occupancyProfileParameters` is an `occupancyProfileParameters` object. If not supplied, it will be generated *de novo* internally. Once again, we recommend to set the parameters of this object in relationship to real ChIP-seq data. `norm = TRUE` and `method` respectively represent if the ChIP-seq like profile should be normalised and if you wish to use an approximation for ChIP-seq profile or not. `moving_kernel` will use `Rcpp` to approximate and compute peaks, `truncated_kernel` will also approximate peaks but without using `Rcpp`, and `exact` will not approximate peaks. These methods represent different way of computing and/or approximating ChIP-seq peaks. Finally, `peakSignificantThreshold` is a threshold at which peaks will be selected. If you select “moving_kernel” then this threshold is a numeric value describing the peak tail height cut-off value. The default in this case is 0.001. In the case of “truncated_kernel” and “exact”, the threshold represents a distance in base pair from the peak summit at which the peak should be cut. In this case, default is set at 1250 base pairs.

It should be noted that these methods will produce very similar results. And by very similar results, we mean nearly identical.

computeChipProfile

To generate a ChIP-seq like profile:

```
chipProfile <- computeChipProfile(setSequence = eveLocus,
  occupancy = Occupancy,occupancyProfileParameters = OPP)
chipProfile
```

The output of this functions is slightly different as it returns a named list (each element in the list is named after the specific combination of *lambda* and *boundMolecules* used to compute occupancies) containing a `GRangesList` of `GRanges` with ChIP profile values as a metadata column. These `GRanges` also differ in the sense that they now contain the whole loci (or accessible loci) cut into bins of size equal to `stepSize` (See `stepSize` slot in `occupancyProfileParameters`). Each `GRangesList` contains `GRanges` for each Loci of interest.

Searching through SitesAboveThreshold and ChIP-seq profiles

As described previously, The size of the `AllSitesAboveThreshold` slot will increase drastically as the number of values assigned to `ScalingFactorPWM` (or `lambda`) and `boundMolecules` increases. In order to navigate

and search this slot with ease, it is possible to use the `searchSites` function. This function may also be used on predicted ChIP-seq profiles (result of `computeChIPProfile`). `searchSites` comes in the following form:

```
searchSites(Sites,ScalingFactor="all", BoundMolecules="all",Locus="all")
```

It is possible to use this function as a simple extraction method similarly to the `AllSitesAboveThreshold` method. In this case, the usage is the following:

```
searchSites(Occupancy)
```

```
## $`lambda` = 1.5 & boundMolecules = 1000`
## GRangesList object of length 1:
## $eve
## GRanges object with 412 ranges and 2 metadata columns:
##      seqnames      ranges strand |      PWMscore
##      <Rle>        <IRanges> <Rle> |      <numeric>
## eve chr2R [5860705, 5860712] + | -1.51655573585429
## eve chr2R [5860709, 5860716] + | -5.33217184502491
## eve chr2R [5860715, 5860722] + |  9.13992557549757
## eve chr2R [5860728, 5860735] + |  5.05434682102833
## eve chr2R [5860758, 5860765] + | -5.15370980167748
## ...      ...      ...      ...
## eve chr2R [5876629, 5876636] + |  5.60817413411963
## eve chr2R [5876635, 5876642] + |  0.202790199774102
## eve chr2R [5876641, 5876648] - | -4.47385601266488
## eve chr2R [5876666, 5876673] + |  2.21133362723558
## eve chr2R [5876684, 5876691] + | -2.28895797651261
##      Occupancy
##      <numeric>
## eve 0.0138683203024566
## eve 0.0138160293072631
## eve 0.0783704718441574
## eve 0.0183246335750422
## eve 0.0138165926852252
## ...      ...
## eve 0.0203268664876583
## eve 0.0139901018008407
## eve 0.0138194725807681
## eve 0.014492381648295
## eve 0.0138454813805688
##
## -----
## seqinfo: 1 sequence from an unspecified genome; no seqlengths
```

If you wish to navigate and extract only certain combinations of `ScalingFactorPWM` and/or `boundMolecules` and/or `Loci`, `searchSites` could be use as shown below:

```
searchSites(chipProfile, ScalingFactor=c(1.5,2.5), BoundMolecules=c(1000,1500)
, Locus=c("eve", "odd"))
```

```
## $`lambda` = 1.5 & boundMolecules = 1000`
## $`lambda` = 1.5 & boundMolecules = 1000`$eve
## GRanges object with 1600 ranges and 1 metadata column:
##      seqnames      ranges strand |      ChIP
##      <Rle>        <IRanges> <Rle> |      <numeric>
## eve chr2R [5860693, 5860703] * |  0.0467998729244692
## eve chr2R [5860703, 5860713] * |  0.051053053031132
```

```
## eve chr2R [5860713, 5860723] * | 0.0554324704104929
## eve chr2R [5860723, 5860733] * | 0.059949075887137
## eve chr2R [5860733, 5860743] * | 0.0646141633273505
## ... ...
## eve chr2R [5876643, 5876653] * | 0.0158243321666681
## eve chr2R [5876653, 5876663] * | 0.0149496630784295
## eve chr2R [5876663, 5876673] * | 0.0140710281760898
## eve chr2R [5876673, 5876683] * | 0.0131862304147329
## eve chr2R [5876683, 5876693] * | 0.0122930573390848
## -----
## seqinfo: 1 sequence from an unspecified genome; no seqlengths
```

Estimating the accuracy of the model

In order to determine how accurate the predicted model is, it is possible to compare the predicted *ChIP-seq like profile* (as built in `computeChipProfile`) to real ChIP-seq data for a given Transcription Factors at loci of interest. `profileAccuracyEstimate` provides a way to compare both profiles. The arguments for this function are the following:

```
profileAccuracyEstimate(LocusProfile,
  predictedProfile, occupancyProfileParameters = NULL)
```

Input data - Accuracy Estimate

`profileAccuracyEstimate` requires only three arguments. `precitedProfile` is the result of `computeChipProfile` and `occupancyProfileParameters` is a `occupancyProfileParameters`. Finally, `LocusProfile` is a list containing actual ChIP-seq profiles. These profiles should be normalised to a base pair level. In other words, a peak should be divided by its width. We also strongly recommend that each loci in `LocusProfile` (each element of the list) should be named in an identical manner as the loci used in `setSequence` (See previous functions). This list should come in the following format:

```
str(eveLocusChip)
```

```
## List of 1
## $ eve: num [1:16000] 0.00755 0.00755 0.00755 0.00755 0.00755 ...
```

In this example, there is only one element in the list. However, this list can be as long as you wish and contain all the Loci that you are interested in.

profileAccuracyEstimate

To test the accuracy the model against ChIP-seq data:

```
AccuracyEstimate <- profileAccuracyEstimate(LocusProfile = eveLocusChip,
  predictedProfile = chipProfile, occupancyProfileParameters = OPP)
AccuracyEstimate
```

The result of this function will be a list of accuracy estimates for every loci and every combination of `ScalingFactorPWM` and `boundMolecules`. The correlation and Mean Squared Error (MSE) represents the correlation and MSE between the predicted profile (for a given combination on `lambda` and `boundMolecules`) and the ChIP-seq profile for the same loci. `meanCorr` and `meanMSE` describe the average correlation and MSE for all loci (for a given combination on `ScalingFactorPWM` and `boundMolecules`). The idea behind average correlation and MSE is that the scaling factor and number of molecules should be the same regardless of the loci as all TF's are contained within the same nucleus. Finally, `meanTheta` is an in house metric describing a

modified ratio of correlation over MSE. The goal is to find the sweet spot between high correlation and low MSE (see `computeOptimal` and `plotOptimalHeatMaps`).

Finding optimal Parameters

As described previously, it is not always possible to know the optimal set of parameters for `ScalingFactorPWM` and `boundMolecules`. `ChIPAnalyser` offers the possibility to backward infer the parameters using the `computeOptimal` function. By testing different combinations of `ScalingFactorPWM` and `boundMolecules`, this function will return the combination with the highest correlation, lowest Mean Squared Error or highest theta depending on which parameter was selected. As a reminder, theta is an in house metric representing a modified ratio of correlation over MSE (extreme values are replaced by threshold). The goal is to find the sweet spot between high correlation and low MSE. Values that should be tested for `ScalingFactorPWM` and for `boundMolecules` should be provided by user. If these values are not provided (default value and only one value for each parameter), then they will be assigned internally. The internal values are the following:

```
ScalingFactorPWM(genomicProfileParameters) <- c(0.25, 0.5, 0.75, 1, 1.25,
1.5, 1.75, 2, 2.5, 3, 3.5, 4, 4.5, 5)

boundMolecules(occupancyProfileParameters) <- c(1, 10, 20, 50, 100,
200, 500, 1000, 2000, 5000, 10000, 20000, 50000, 100000,
200000, 500000, 1000000)
```

In terms of its arguments, `computeOptimal` can be described as:

```
computeOptimal(DNASequenceSet,
  genomicProfileParameters,
  LocusProfile,
  setSequence,
  DNAAccessibility = NULL,
  occupancyProfileParameters = NULL,
  parameter = "all",
  peakMethod="moving_kernel")
```

Please note that this functions will take some time to complete. Do not be alarmed if it seems to have stalled.

Input Data - Optimal Parameters

`computeOptimal` is essentially a combination of previous functions (with a bit more magic to it). For this reason, data input is extremely similar to the functions described above. As a quick reminder:

- `DNASequenceSet`, a `DNAStringSet` (or `BSgenome`) containing the sequences of the organism of interest.
- `genomicProfileParameters`, a `genomicProfileParameters` object containing at least a *Position Weight Matrix* or *Position Frequency Matrix*. All other slots will be computed internally.
- `LocusProfile`, a named list of ChIP-seq profile for loci of interest.
- `setSequence`, a named `GRanges` containing loci of interest.
- `DNAAccessibility`, a `GRanges` containing Accessible DNA.
- `occupancyProfileParameters`, an `occupancyProfileParameters` object. Although optional, we strongly advise to tailor this object by using values directly extracted from `LocusProfile`

`parameter` defines which metric you wish to compute. There are four possible choices: *correlation*, *MSE*, *theta* or *all*. It is imperative that the lists/`GRanges` are named with the name of the Loci of interest. `peakMethod` describes if you wish to use an approximation for ChIP-seq profile peaks. `moving_kernel` will use `Rcpp` to approximate and compute peaks, `truncated_kernel` will also approximate peaks but without

using Rcpp, and `exact` will not approximate peaks. These methods represent different way of computing and/or approximating ChIP-seq peaks.

computeOptimal

As a example describing the usage of `compute optimal`

```
optimalParam <- computeOptimal(DNASequenceSet = DNASequenceSet,
  genomicProfileParameters = GPP,
  LocusProfile = eveLocusChip,
  setSequence = eveLocus,
  DNAAccessibility = Access,
  occupancyProfileParameters = OPP,
  parameter = "all")
optimalParam
```

This functions returns either a list or a list of lists (if “all” parameter was selected). Each element in the list represents the **optimal set of parameters**, the **optimal matrix** (a matrix with correlation, MSE and/or theta computed for a given combination of `ScalingFactorPWM` and `boundMolecules`) and finally the **selected parameter**.

Plotting Results

As it is the case in many fields, data visualisation is a key aspect in any analysis. For this purpose, ChIPAnalyser offers two plotting functions: `plotOptimalHeatMaps` and `plotOccupancyProfile`.

Optimal Parameters

Once you have computed the optimal set of parameters, it is possible to plot these results in the form of a heat map using `plotOptimalHeatMaps`. Depending on what you are interested in, this function will either plot *correlation*, *MSE*, *theta* or *all of the previous*. This functions requires minimal input as described below:

```
plotOptimalHeatMaps(optimalParam=optimalParam ,
  parameter="all", Contour=TRUE)
```

Input Data & Plotting

`plotOptimalHeatMaps` only requires one data input in the form of the result of `computeOptimal` (see `computeOptimal`). The `parameter` argument defines which of the following parameters you wish to plot: *correlation*, *MSE*, *theta* or *all of the previous*. Finally, `Contour` defines if you which to plot Contour lines on your heat map. As an example:

```
plotOptimalHeatMaps(optimalParam, parameter="all")
```

See plot in **Quick Guide**

The boxed tile represents the highest correlation or theta for a given combination of `ScalingFactorPWM` and `boundMolecules`. In the case of MSE the boxed tile represents the lowest Mean Squared Error.

Plotting Profiles

ChIPAnalyser produces ChIP-seq like profiles. It is possible to plot these profiles but also to add a variety of features to these plots. `plotOccupancyProfile` takes care of plotting with the following arguments:

```
plotOccupancyProfile <- function(predictedProfile,
  setSequence,
  profileAccuracy = NULL,
  chipProfile = NULL,
  occupancy = NULL,
  PWM=FALSE,
  DNAAccessibility = NULL,
  occupancyProfileParameters = NULL,
  geneRef = NULL)
```

Input Data & Profiles

In order to increase plotting flexibility, `plotOccupancyProfile` only plots one profile at a time. In practice, this means that only simple data units should be parsed to this functions. This also means that the main title is left to the user discretion. The arguments described above should come in the following format:

- `precitedProfile`, a GRanges object containing the predicted ChIP-seq like profile for one locus and one combination of `lambda` and `boundMolecules`.
- `setSequence`, a GRanges object containing the locus of interest.
- `profileAccuracy`, the profile Accuracy estimate for one loci and for one combination of `lambda` and `boundMolecules`
- `chipProfile`, a vector containing ChIP-seq data for locus of interest. In previous functions, ChIP-seq data was stored in a named list. In this case, it is the individual numeric vector contained within that list.
- `occupancy`, a GRanges object containing both PWMScore and Occupancy. This GRanges is the result of `computeOccupancy` and should only contain a GRanges object for one locus and one combination of `lambda` and `boundMolecules`.
- `PWM`, a logical operator indicating wherever you wish to plot *occupancy* or *PWMScores*. It is necessary to also include *occupancy* data.
- `DNAAccessibility`, a GRanges object containing DNAAccessibility. DNAAccessibility is similar to DNAAccessibility data described previously.
- `occupancyProfileParameters`, an `occupancyProfileParameters` object. This object should be the same as the one used in functions described above. However, the minimal requirement is that the `stepSize` slot remains consistent with `stepSize` used previously. As a reminder, `stepSize` default value is set at 10.
- `geneRef`, a List containing genetic information (3'UTR, 5'UTR, exons, intron and enhancers). Each element of this list, is a GRanges containing the information regarding 3'UTR, 5'UTR, exons, intron and enhancers.

As this object has not yet be described, `geneRef` should come in a similar format as the following:

```
geneRef
```

```
## $exon
## GRanges object with 26713 ranges and 0 metadata columns:
##           seqnames      ranges strand
##           <Rle>         <IRanges> <Rle>
##   CG17683   chr2R      [18442, 18629]   +
##   CG17683   chr2R      [18681, 18773]   +
##   CG17683   chr2R      [18827, 19484]   +
```

```

## CG17683 chr2R [19542, 20468] +
## CG17683 chr2R [18442, 18629] +
## ... ...
## CG33680 chr2R [21137781, 21137839] -
## CG30428 chr2R [21140837, 21140963] +
## CG30428 chr2R [21141104, 21141284] +
## CG30428 chr2R [21141343, 21141601] +
## CG30428 chr2R [21141651, 21142371] +
## -----
## seqinfo: 14 sequences from an unspecified genome; no seqlengths
##
## $intron
## GRanges object with 22058 ranges and 0 metadata columns:
##      seqnames      ranges strand
##      <Rle>         <IRanges> <Rle>
## CG17683 chr2R [18630, 18680] +
## CG17683 chr2R [18774, 18826] +
## CG17683 chr2R [19485, 19541] +
## CG17683 chr2R [18630, 18692] +
## CG17683 chr2R [18774, 18826] +
## ... ...
## CG33680 chr2R [21137114, 21137174] -
## CG33680 chr2R [21137423, 21137780] -
## CG30428 chr2R [21140964, 21141103] +
## CG30428 chr2R [21141285, 21141342] +
## CG30428 chr2R [21141602, 21141650] +
## -----
## seqinfo: 13 sequences from an unspecified genome; no seqlengths
##
## $`5UTR`
## GRanges object with 6029 ranges and 0 metadata columns:
##      seqnames      ranges strand
##      <Rle>         <IRanges> <Rle>
## CG17683 chr2R [18442, 18566] +
## CG17683 chr2R [18442, 18566] +
## CG17683 chr2R [18487, 18629] +
## CG17683 chr2R [18681, 18811] +
## CG17683 chr2R [18498, 18773] +
## ... ...
## CG9380 chr2R [21076340, 21076360] -
## CG9380 chr2R [21076340, 21076360] -
## Kr chr2R [21114138, 21114474] +
## CG30429 chr2R [21133990, 21134051] +
## CG30428 chr2R [21140837, 21140961] +
## -----
## seqinfo: 13 sequences from an unspecified genome; no seqlengths
##
## $`3UTR`
## GRanges object with 4556 ranges and 0 metadata columns:
##      seqnames      ranges strand
##      <Rle>         <IRanges> <Rle>
## CG17683 chr2R [20162, 20468] +
## CG17683 chr2R [20162, 20468] +
## CG17683 chr2R [20162, 20468] +

```

```
## CG17683 chr2R [20162, 20468] +
## CG17683 chr2R [20162, 20468] +
## ... ...
## CG9380 chr2R [21072649, 21072809] -
## Kr chr2R [21116357, 21117057] +
## CG30429 chr2R [21135028, 21135109] +
## CG33680 chr2R [21136529, 21136529] -
## CG30428 chr2R [21142001, 21142371] +
## -----
## seqinfo: 13 sequences from an unspecified genome; no seqlengths
```

It should be noted that only two arguments are necessary (`predictedProfile` and `setSequence`). The more arguments are provided the more information will be plotted. As an example:

```
plotOccupancyProfile(predictedProfile=chipProfile[[1]][[1]],
  setSequence=eveLocus,
  profileAccuracy = AccuracyEstimate[[1]][[1]],
  chipProfile = eveLocusChip[[1]],
  occupancy = AllSitesAboveThreshold(Occupancy)[[1]][[1]],
  DNAAccessibility = Access,
  occupancyProfileParameters = OPP,
  geneRef =geneRef)
```

Session Information

```
sessionInfo()
```

```
## R version 3.3.2 (2016-10-31)
## Platform: x86_64-pc-linux-gnu (64-bit)
## Running under: Ubuntu 16.04.3 LTS
##
## locale:
##  [1] LC_CTYPE=en_GB.UTF-8      LC_NUMERIC=C
##  [3] LC_TIME=en_GB.UTF-8      LC_COLLATE=en_GB.UTF-8
##  [5] LC_MONETARY=en_GB.UTF-8  LC_MESSAGES=en_GB.UTF-8
##  [7] LC_PAPER=en_GB.UTF-8     LC_NAME=C
##  [9] LC_ADDRESS=C             LC_TELEPHONE=C
## [11] LC_MEASUREMENT=en_GB.UTF-8 LC_IDENTIFICATION=C
##
## attached base packages:
## [1] parallel stats4      stats      graphics  grDevices  utils      datasets
## [8] methods   base
##
## other attached packages:
## [1] BSgenome.Dmelanogaster.UCSC.dm3_1.4.0
## [2] ChIPanalyser_0.99.12
## [3] RcppRoll_0.2.2
## [4] BSgenome_1.42.0
## [5] rtracklayer_1.34.2
## [6] Biostrings_2.42.1
## [7] XVector_0.14.1
## [8] GenomicRanges_1.26.4
## [9] GenomeInfoDb_1.10.3
```

```
## [10] IRanges_2.8.2
## [11] S4Vectors_0.12.2
## [12] BiocGenerics_0.20.0
##
## loaded via a namespace (and not attached):
## [1] Rcpp_0.12.12      knitr_1.17
## [3] magrittr_1.5      GenomicAlignments_1.10.1
## [5] zlibbioc_1.20.0   BiocParallel_1.8.2
## [7] lattice_0.20-35   stringr_1.2.0
## [9] tools_3.3.2       grid_3.3.2
## [11] SummarizedExperiment_1.4.0 Biobase_2.34.0
## [13] htmltools_0.3.6   yaml_2.1.14
## [15] rprojroot_1.2     digest_0.6.12
## [17] Matrix_1.2-10     bitops_1.0-6
## [19] RCurl_1.95-4.8    evaluate_0.10.1
## [21] rmarkdown_1.6     stringi_1.1.5
## [23] backports_1.1.0   Rsamtools_1.26.2
## [25] XML_3.98-1.9
```

References

Zabet NR, Adryan B (2015) Estimating binding properties of transcription factors from genome-wide binding profiles. *Nucleic Acids Res.*, 43, 84–94.