

methylPipe

Kamal Kishore, Mattia Pelizzola

May 19, 2021

Contents

1	Introduction	1
2	Quick overview	1
3	Profiling Genome wide DNA methylation	2
4	Descriptive statistics of DNA methylation	5
5	Profiling DNA methylation on genomic ranges	6
6	Differential DNA methylation	8
7	Session Information	11

1 Introduction

In this document you can find a brief tutorial on the *methylPipe* package for the analysis of base-pair resolution DNA methylation data. DNA methylation is a potentially heritable epigenetic modification of the genomic DNA typical of most eukaryotic organisms and critical for the regulation of gene transcription. The level of DNA methylation varies according to age, diet and environment and an appropriate control of methylation patterns is important for the onset of cellular differentiation processes and can be deregulated in diseases as cancer. In eukaryotic genomes the methyl group can be added to Cytosines in the CG, CHG and CHH sequence context (H being one of A, C or T). It is nowadays possible to generate genome-wide base-pair resolution DNA methylation maps (see Lister R et al Nature 2009). *methylPipe* is an R package that includes a series of objects and methods for a memory efficient management, query, analysis and visualization of DNA methylation data and their integration with heterogeneous data types. This package has been developed alongwith *compEpiTools* which provide functions and methods for the analysis of epigenomics data. The data package *ListerEtAlBSseq* has also been developed which consists of base resolution Bisulfite sequencing data of H1 and IMR90 cell line (Lister R et al Nature 2009).

The overview section of this document will briefly walk you through the main available functionalities, while the following sections will provide additional details on a selection of available classes and methods.

2 Quick overview

A number of **classes** are defined in *methylPipe*: *BSdata*, *BSdataSet*, *GEcollection* and *GElis*.

- *BSdata* is a reference class to store base resolution DNA methylation data generated from a high-throughput sequencing experiment for a given biological sample.

- The *BSdataSet* class allows combining DNA methylation data for several samples for the same organism.
- The *GEcollection* class is used to store the DNA methylation status of a collection of genomic regions.
- The *GElist* class is the list of multiple objects of class *GEcollection*.

methylPipe consists of methods which allows analysis and visualisation of genome wide DNA methylation data. The `meth.call` function processes methylation information from aligned files and creates tabular data file. `BSprepare` processes and tabix compresses such tabular data files for creation of a *BSdata* object for each sample. The *BSdataSet* classe stores multiple such samples of class *BSdata*. The `mCsmoothing` methods extracts DNA methylation data for a given genomic region and smooths/plots it.

For methylation analysis on your regions of interest *GEcollection* and *GElist* acts as a repository of methylation data across various genomic ranges. `getCpos` and `getCposDensity` retrieves all potential genomic cytosine positions and their density across genomic regions, repectively. `mapBSdata2GRanges` allows extraction of DNA methylation information across multiple genomic ranges and `profileDNAMetBin` determines the absolute and relative methylation to create a *GEcollection* object. Different methods for the detection of differentially methylated regions are implemented, according to the number of samples (see `findDMR` and `consolidateDMRs`). In addition the methylation profile can be visualised using `plotMeth`. Moreover, integrative analysis with other epignomics dataset can be performed by using `heatmapdata`, `heatmapPlot` functions of the package *compEpiTools*.

The methods that are most demanding in terms of computational resources are optimized for low memory fingerprint and multi-processor support. *methylPipe* includes a series of objects and methods that can be used as building blocks for the creation of pipelines for the data analysis of epigenomics data, with particular emphasis on DNA methylation, and their integration with any kind of annotation or additional data type.

3 Profiling Genome wide DNA methylation

First of all the *methylPipe* and the genome sequence libraries are loaded:

```
library(methylPipe)
library(BSgenome.Hsapiens.UCSC.hg18)
```

The DNA methylation information can be read from a text file. See the documentation of the `BSprepare` on how to build data suitable to populate a *BSdata* object. Moreover, the methylation information can also be read directly from the sorted SAM file(s) generated from the BISMARCK aligner. The function `meth.call` process the sorted SAM file. The user can specify the sequence context in which the methylation information is read from these files either "CpG" or "All". In case of whole genome Bisulfite data especially in case of non-CG methylation the number of potential methylation sites are enormous and vast majority of them are unmethylated. In order to avoid storing too much data while maintaining the ability to identify methylated, unmethylated and uncovered Cytosines, *methylPipe* only stores and access C positions with at least 1 mC read. Genomic regions not covered by sequencing are stored as a *GRanges* object. The function `meth.call` produces methylation call text file and uncovered genomic regions file for each sample in the output folder. Finally, when profiling DNA methylation unmethylated Cs are determined as those genomic cytosines not having any methylated reads and not belonging to uncovered genomic regions.

```
file_loc <- system.file('extdata', 'test_methcall', package='methylPipe')
meth.call(files_location=file_loc, output_folder=tempdir(), no_overlap=TRUE,
          read.context="CpG", Nproc=1)
```

methylPipe adopts TABIX compressed indexing of flatfiles to reduce disk space which allows fast access. The *methylPipe* library includes a small subset of the first two base resolution human DNA methylomes (Lister R et al Nature 2009) for two well known human cell lines: embryonic stem cells (H1) and fetal lung fibroblasts (IMR90). The methylation data can be stored into *BSdata* objects and then collected together within a *BSdataSet* object. The *BSprepare* command is not run in the following example since it requires a local copy of tabix software. Hence, for this example we use pre-generated tabix indexed files from *methylPipe* for H1 and IMR90.

```
file_loc <- system.file('extdata', 'H1_chr20_CG_10k.txt', package='methylPipe')
#BSprepare(files_location=file_loc, output_folder=file_loc, tabix="/path-to-tabix/")
uncov_GR <- GRanges(Rle('chr20'), IRanges(c(14350,69251,84185), c(18349,73250,88184)))
H1data <- system.file('extdata', 'H1_chr20_CG_10k_tabix_out.txt.gz',
                    package='methylPipe')
H1.db <- BSdata(file=H1data, uncov=uncov_GR, org=Hsapiens)
H1.db

## GRanges object with 3 ranges and 0 metadata columns:
##      seqnames      ranges strand
##      <Rle>      <IRanges> <Rle>
## [1]   chr20 14350-18349      *
## [2]   chr20 69251-73250      *
## [3]   chr20 84185-88184      *
## -----
## seqinfo: 1 sequence from an unspecified genome; no seqlengths
## [1] "chr20"
## GRanges object with 6 ranges and 4 metadata columns:
##      seqnames      ranges strand | Context      C      T
##      <Rle> <IRanges> <Rle> | <character> <numeric> <numeric>
## [1]   chr20      8179      + |      CG          2          4
## [2]   chr20      8180      - |      CG          4          4
## [3]   chr20      8426      + |      CG          1          0
## [4]   chr20      8427      - |      CG          5          0
## [5]   chr20      8432      + |      CG          1          0
## [6]   chr20      8433      - |      CG          6          0
##      Significance
##      <numeric>
## [1]          20
## [2]          48
## [3]          14
## [4]          84
## [5]          14
## [6]         102
## -----
## seqinfo: 1 sequence from an unspecified genome; no seqlengths
```

Multiple *BSdata* objects can be stored in *BSdataSet* object by specifying group name for each sample either as "C" (control) or "E" (Experiment):

```

IMR90data <- system.file('extdata', 'IMR90_chr20_CG_10k_tabix_out.txt.gz',
                        package='methylPipe')
IMR90.db <- BSdata(file=IMR90data, uncov=uncov_GR, org=Hsapiens)
H1.IMR90.set <- BSdataSet(org=Hsapiens, group=c("C","E"), IMR90=IMR90.db, H1=H1.db)
H1.IMR90.set

## S4 Object of class BSdataSet
##
## BSdata objects contained:
## [1] "IMR90" "H1"
##
## Groups of objects:
## [1] "C" "E"
##
## Associated organism genome:
## Homo sapiens
##

```

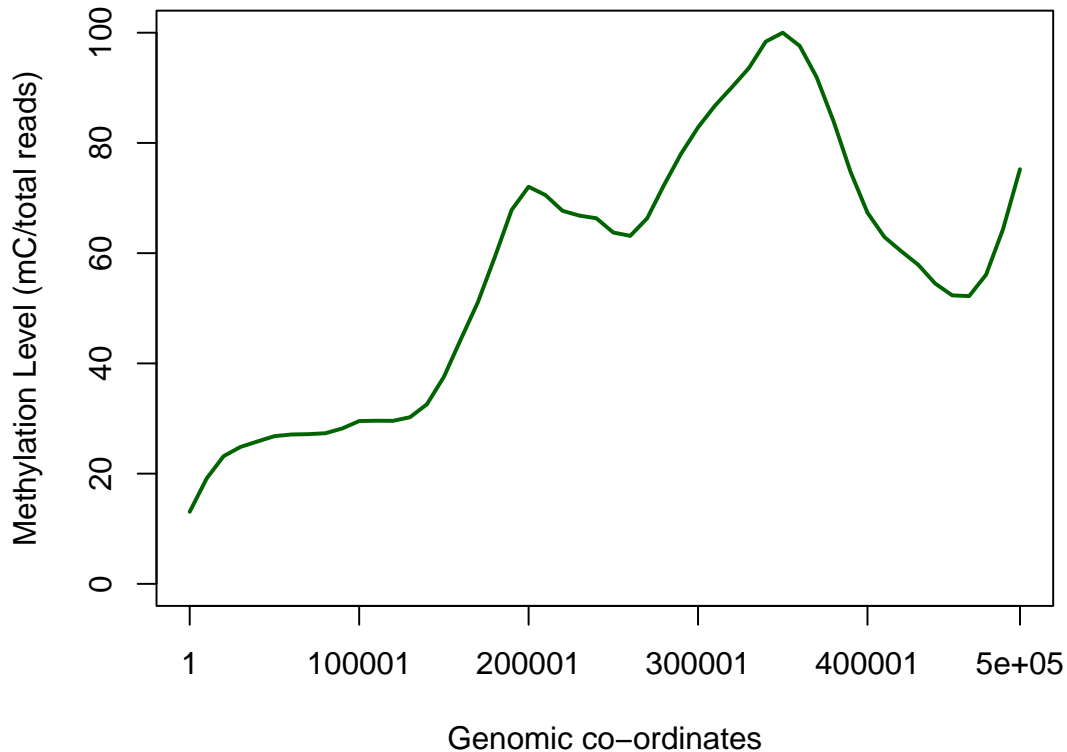
The data includes the chromosome, the genomic position, strand and sequence context for each methyl-Cytosine (mC). In addition, for each mC the number of reads with C at that genomic position and strand (supporting DNA methylation), the number of reads with T (evidence of an unmethylated C) and the binomial p-value supporting the mC call are attached. Importantly, these data are referenced into the *BSdata* but they are not actually loaded into the memory. The `mCsmoothing` method allow uploading into memory the DNA methylation data for a given genomic region and to display their methylation profile over it. The smoothing can be performed by either "sum" or "mean" of methylation level for each regions.

```

gr <- GRanges("chr20", IRanges(1, 5e5))
sres <- mCsmoothing(H1.db, gr, Scorefun='sum', Nbins=50, plot=TRUE)

```

Smoothed methylation levels



4 Descriptive statistics of DNA methylation

methyPipe allows checking the basic stats about the methylation data such as range, mean and quantile distribution of methylation and assess sample similarity with correlation and clustering analysis. The `methstats` method computes pairwise correlation coefficients (Pearson) between the methylation profiles across all the samples in *BSdataSet* object. It outputs scatter plot matrix of correlation coefficients. Finally, it performs (euclidean distance based) hierarchical clustering of samples and outputs the dendrogram. In the example below, the analysis is performed on *BSdataSet* object of artificially replicated H1 and IMR90.

```
stats.set <- BSdataSet(org=Hsapiens, group=c("C","C","E","E"), IMR_1=IMR90.db,
IMR_2=IMR90.db, H1_1=H1.db,H1_2=H1.db)
stats_res <- methstats(stats.set,chrom="chr20",mcClass='mCG', Nproc=1)
stats_res
```

```
## $descriptive_stats
##      IMR_1      IMR_2      H1_1      H1_2
##  Min.   :0.0000  Min.   :0.0000  Min.   :0.0000  Min.   :0.0000
## 1st Qu.:0.5360  1st Qu.:0.5360  1st Qu.:0.6670  1st Qu.:0.6670
## Median :0.8570  Median :0.8570  Median :0.8480  Median :0.8480
## Mean   :0.7136  Mean   :0.7136  Mean   :0.7497  Mean   :0.7497
## 3rd Qu.:1.0000  3rd Qu.:1.0000  3rd Qu.:0.9700  3rd Qu.:0.9700
```

```
## Max.      :1.0000   Max.      :1.0000   Max.      :1.0000   Max.      :1.0000
##
## $correlation_mat
##           IMR_1      IMR_2      H1_1      H1_2
## IMR_1  1.0000000  1.0000000  0.1653737  0.1653737
## IMR_2  1.0000000  1.0000000  0.1653737  0.1653737
## H1_1   0.1653737  0.1653737  1.0000000  1.0000000
## H1_2   0.1653737  0.1653737  1.0000000  1.0000000
```

5 Profiling DNA methylation on genomic ranges

methylPipe allows profiling DNA methylation over a set of genomic regions and determining their absolute or relative methylation profiles. Absolute methylation is here defined as the genomic density of mC calls per bp, while relative methylation is the proportion of density of methylated genomic sites over the density of potential methylation sites.

The `mapBSdata2GRanges` method retrieves mC calls given a *BSdata* object of a sample and a set of genomic regions. In the following example we will extract the CG sites for a set of genomic regions. In particular we will consider a *GRanges* object of genomic regions that spans 2Kb upstream and downstream the TSS of the transcript ids on the first 500Kb of chromosome 20:

```
gr_file <- system.file('extdata', 'GR_chr20.Rdata', package='methylPipe')
load(gr_file)
resmC <- mapBSdata2GRanges(GenoRanges=GR_chr20, Sample=H1.db, context='CG')
head(resmC[[4]])

## GRanges object with 6 ranges and 4 metadata columns:
##      seqnames      ranges strand |      Context          C          T
##      <Rle> <IRanges> <Rle> | <character> <numeric> <numeric>
## [1] chr20      116504      + |      CG          18          25
## [2] chr20      116757      + |      CG          15           4
## [3] chr20      116758      - |      CG          24          12
## [4] chr20      116796      + |      CG          32           0
## [5] chr20      116797      - |      CG           8           1
## [6] chr20      117132      - |      CG           9           3
##      Significance
##      <numeric>
## [1]          202
## [2]          227
## [3]          334
## [4]          564
## [5]          128
## [6]          132
## -----
## seqinfo: 1 sequence from an unspecified genome; no seqlengths
```

To profile DNA methylation on genomic regions of interest the *GEcollection* class (collection of genomic regions) is used. These genomic regions can be any regions of interest such as one or more gene loci, promoter regions, a set of transcription factor binding sites, a set of regions enriched by some histone mark or a set of transposable elements. Most of these regions can be usually retrieved from either

publications or public databases like the UCSC Genome Browser and its Table Browser tools (see *compEpiTools* for methods for generation/manipulation of these regions). The *GEcollection* class has data slots and methods tailored for storing and analyzing DNA methylation data. It is an extension of the *RangedSummarizedExperiment* class from the *SummarizedExperiment* package.

For the creation of the *GEcollection*, the `getCpos` and `mapBSdata2GRanges` methods are called automatically to populate its various data slots. Further arguments of the latter function allows to filter the mC sites based on the depth of sequencing, number of reads with C at that genomic position and their binomial p-value. The `profileDNAmetBin` method acts as a wrapper for these methods to profile absolute and relative DNA methylation for a set of regions within a *GEcollection*. In particular for each genomic region, and possibly each bin of it, the absolute methylation density (mC/bp), the density of possible methylation sites (C/bp) and the relative methylation level (mC/C) will be determined. For example, in case of the CG context: mCG/bp, CG/bp and mCG/CG densities will be stored for each bin of each genomic region in the `binmC`, `binC` and `binrC` data slots, respectively.

```

gec.H1 <- profileDNAmetBin(GenoRanges=GR_chr20, Sample=H1.db, mcCLASS='mCG', nbins=3)

## Warning in min(assays(object)[["binscore"]], na.rm = TRUE): no non-missing arguments to
## min; returning Inf

binmC(gec.H1)[4:5,]

##           1      2      3
## [1,] 0.00378 0.00911 0.01440
## [2,] 0.00969 0.01170 0.00825

binC(gec.H1)[4:5,]

##           1      2      3
## [1,] 0.0075 0.0165 0.0210
## [2,] 0.0150 0.0180 0.0105

binrC(gec.H1)[4:5,]

##           1      2      3
## [1,] 50.4 55.2 68.8
## [2,] 64.6 64.8 78.6

```

GEcollection objects can be **subsetting and combined**. `Subset` can be useful to work only on the regions on a particular strand and/or chromosome or chromosomal region.

```

gec1 <- gec.H1[start(gec.H1) < 153924]
gec2 <- gec.H1[start(gec.H1) > 153924]

```

Multiple *GEcollection* objects can be saved into a *GElist* object. This can be convenient if there are many of them and if one would like to iteratively apply the same method to them, for example profiling DNA methylation for their genomic regions in the same sample.

```

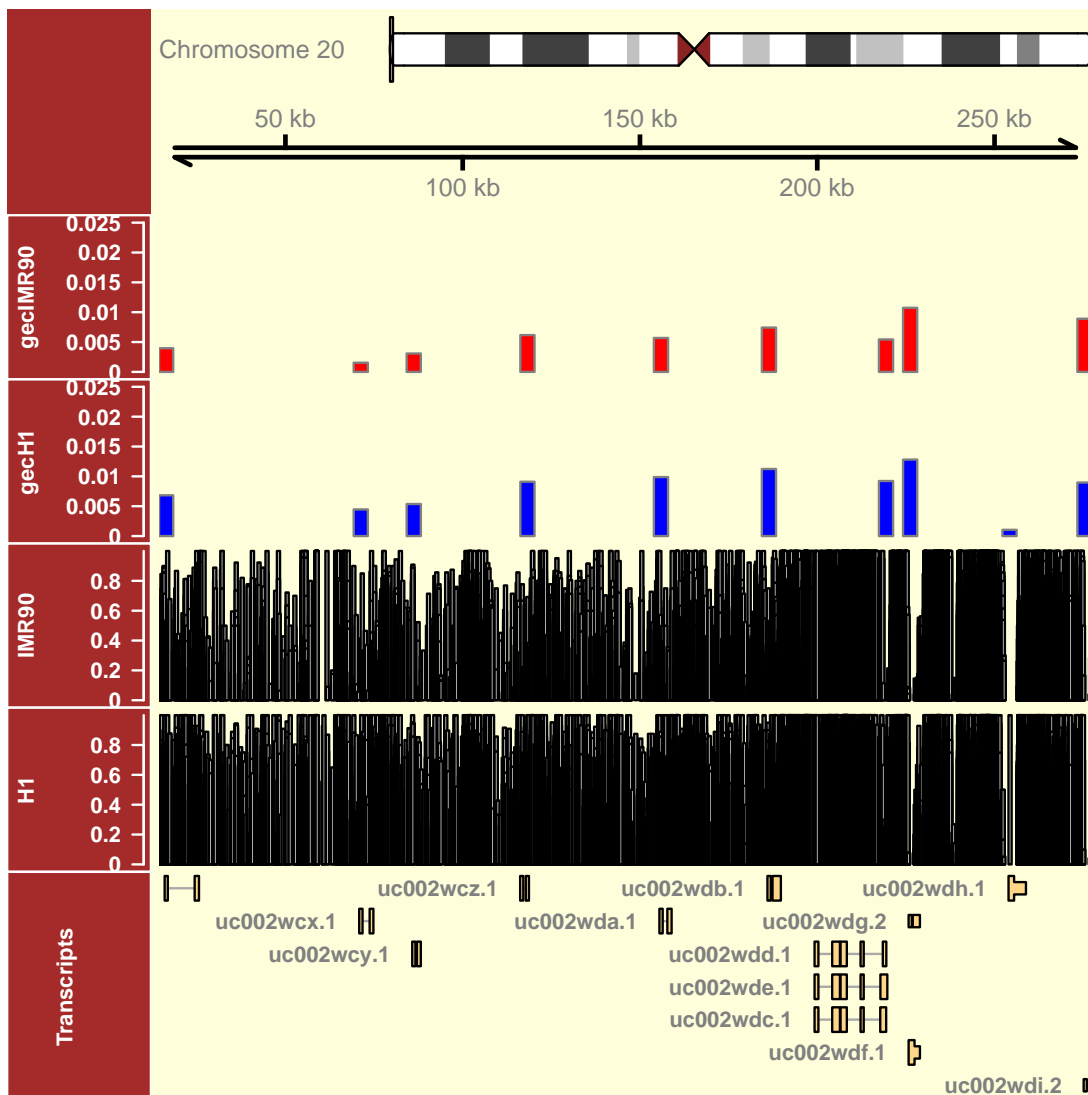
gecIMR_file <- system.file('extdata', 'gec.IMR90.Rdata', package='methylPipe')
load(gecIMR_file)
gel <- GElist(gecIMR90=gec.IMR90, gecH1=gec.H1)
print(names(gel))

```

```
## [1] "gecIMR90" "gecH1"
```

The *GElist* objects can be visualized using `plotMeth`. This allows methylation data of various samples to be displayed together with the annotation information for the genomic regions of interest. Moreover, various epigenomics data tracks can also be visualized together with the methylation information. See the documentation of the `plotMeth` for more details.

```
library(TxDb.Hsapiens.UCSC.hg18.knownGene)
txdb <- TxDb.Hsapiens.UCSC.hg18.knownGene
gel <- GElist(gecIMR90=gec.IMR90[1:10], gecH1=gec.H1[1:10])
plotMeth(gel, colors=c("red","blue"), datatype=c("mC","mC"), yLim=c(.025, .025), brmeth=list(IMR90,
  mcContext="CG", transcriptDB=txdb, chr="chr20", start=14350, end=277370, org=Hsapiens)
```



6 Differential DNA methylation

The `findDMR` function can be used to identify differentially methylated regions (DMRs) for any sequence context (mCG, mCHG or mCHH). DMRs can be identified comparing the mC methylation levels (the

proportion of methylated over total reads for a set of mC) between **two samples** or within a group of **N samples**, with the Wilcoxon or the Kruskal-Wallis non parametric test, respectively. Only cytosine positions covered in both the samples (pairwise analysis) or all the samples (multiple sample analysis) are considered for DMR identification. The algorithm uses dynamic sliding window approach which identifies DMRs depending on methylated cytosine frequency and their relative distance. Moreover, to perform regions specific analysis (such as promoter or CpG island centric) genomic regions could be specified and the algorithm will report DMRs only within those regions, if any. In a first step, all the evaluated regions are reported in the output with their corresponding statistical significance. The methylation difference is determined in terms of MethDiffPerc (percentage difference between the mean methylation of experiment and control) and log2Enrichment (log2 of mean methylation of experiment over control).

```
DMRs <- findDMR(object= H1.IMR90.set, Nproc=1, ROI=GR_chr20, MCClass='mCG',
  dmrSize=6, dmrBp=800)
head(DMRs)

## GRanges object with 6 ranges and 3 metadata columns:
##      seqnames      ranges strand |      pValue MethDiff_Perc
##      <Rle>      <IRanges> <Rle> | <numeric>      <numeric>
## [1]   chr20 14404-15060      * |      0.036          37.133
## [2]   chr20 15001-15095      * |      0.036          41.817
## [3]   chr20 15760-15906      * |      0.787           8.800
## [4]   chr20 16170-16792      * |      0.059          40.033
## [5]   chr20 16721-17355      * |      0.059          35.133
## [6]   chr20 69287-69975      * |      0.059          40.033
##      log2Enrichment
##      <numeric>
## [1]          1.486
## [2]          1.341
## [3]          0.181
## [4]          1.188
## [5]          0.917
## [6]          1.277
## -----
## seqinfo: 1 sequence from an unspecified genome; no seqlengths
```

The `consolidateDMRs` function can be used to multiple-testing correct the DMRs and consolidate them according to their relative distance, type of DMRs and thresholds of p-value/Methylation difference/log enrichment. A final *GRanges* object with the set of DMRs including p-value, methylation difference and log enrichment is provided:

```
hyper.DMRs.conso <- consolidateDMRs(DmrGR=DMRs, pvThr=0.05, GAP=100, type="hyper")
hyper.DMRs.conso[1:4]

## GRanges object with 4 ranges and 3 metadata columns:
##      seqnames      ranges strand |      pValue MethDiff_Perc
##      <Rle>      <IRanges> <Rle> | <numeric>      <numeric>
## [1]   chr20 14404-15095      * |      0.010          39.475
## [2]   chr20 69927-70629      * |      0.036          32.717
## [3]   chr20 72059-72852      * |      0.010          52.458
## [4]   chr20 84742-85118      * |      0.036          36.750
```

```
##      log2Enrichment
##      <numeric>
## [1]      1.414
## [2]      1.320
## [3]      2.037
## [4]      0.921
## -----
## seqinfo: 1 sequence from an unspecified genome; no seqlengths
```

7 Session Information

```
sessionInfo()

## R version 4.1.0 beta (2021-05-03 r80259)
## Platform: x86_64-pc-linux-gnu (64-bit)
## Running under: Ubuntu 20.04.2 LTS
##
## Matrix products: default
## BLAS: /home/biocbuild/bbs-3.14-bioc/R/lib/libRblas.so
## LAPACK: /home/biocbuild/bbs-3.14-bioc/R/lib/libRlapack.so
##
## locale:
## [1] LC_CTYPE=en_US.UTF-8      LC_NUMERIC=C
## [3] LC_TIME=en_GB              LC_COLLATE=C
## [5] LC_MONETARY=en_US.UTF-8    LC_MESSAGES=en_US.UTF-8
## [7] LC_PAPER=en_US.UTF-8      LC_NAME=C
## [9] LC_ADDRESS=C               LC_TELEPHONE=C
## [11] LC_MEASUREMENT=en_US.UTF-8 LC_IDENTIFICATION=C
##
## attached base packages:
## [1] parallel stats4 stats graphics grDevices utils
## [7] datasets methods base
##
## other attached packages:
## [1] TxDb.Hsapiens.UCSC.hg18.knownGene_3.2.2
## [2] GenomicFeatures_1.45.0
## [3] AnnotationDbi_1.55.0
## [4] BSgenome.Hsapiens.UCSC.hg18_1.3.1000
## [5] BSgenome_1.61.0
## [6] rtracklayer_1.53.0
## [7] methylPipe_1.27.0
## [8] Rsamtools_2.9.0
## [9] Biostrings_2.61.0
## [10] XVector_0.33.0
## [11] SummarizedExperiment_1.23.0
## [12] Biobase_2.53.0
## [13] MatrixGenerics_1.5.0
## [14] matrixStats_0.58.0
## [15] GenomicRanges_1.45.0
## [16] GenomeInfoDb_1.29.0
## [17] IRanges_2.27.0
## [18] S4Vectors_0.31.0
## [19] BiocGenerics_0.39.0
##
## loaded via a namespace (and not attached):
## [1] colorspace_2.0-1      rjson_0.2.20
## [3] ellipsis_0.3.2        biovizBase_1.41.0
## [5] htmlTable_2.2.1       base64enc_0.1-3
```

```

## [7] dichromat_2.0-0      rstudioapi_0.13
## [9] bit64_4.0.5          fansi_0.4.2
## [11] codetools_0.2-18     splines_4.1.0
## [13] cachem_1.0.5         knitr_1.33
## [15] Formula_1.2-4       cluster_2.1.2
## [17] dbplyr_2.1.1        png_0.1-7
## [19] compiler_4.1.0      httr_1.4.2
## [21] backports_1.2.1     assertthat_0.2.1
## [23] Matrix_1.3-3        fastmap_1.1.0
## [25] lazyeval_0.2.2      limma_3.49.0
## [27] htmltools_0.5.1.1   prettyunits_1.1.1
## [29] tools_4.1.0         gtable_0.3.0
## [31] glue_1.4.2          GenomeInfoDbData_1.2.6
## [33] dplyr_1.0.6         rappdirs_0.3.3
## [35] Rcpp_1.0.6          vctrs_0.3.8
## [37] xfun_0.23           stringr_1.4.0
## [39] lifecycle_1.0.0     restfulr_0.0.13
## [41] ensemblDb_2.17.0    gtools_3.8.2
## [43] XML_3.99-0.6        zlibbioc_1.39.0
## [45] scales_1.1.1        VariantAnnotation_1.39.0
## [47] hms_1.1.0           ProtGenerics_1.25.0
## [49] AnnotationFilter_1.17.0 RColorBrewer_1.1-2
## [51] yaml_2.2.1          curl_4.3.1
## [53] memoise_2.0.0       gridExtra_2.3
## [55] ggplot2_3.3.3       biomaRt_2.49.0
## [57] rpart_4.1-15        latticeExtra_0.6-29
## [59] stringi_1.6.2       RSQLite_2.2.7
## [61] highr_0.9           BiocIO_1.3.0
## [63] checkmate_2.0.0     caTools_1.18.2
## [65] filelock_1.0.2      BiocParallel_1.27.0
## [67] rlang_0.4.11        pkgconfig_2.0.3
## [69] bitops_1.0-7        evaluate_0.14
## [71] lattice_0.20-44     purrr_0.3.4
## [73] GenomicAlignments_1.29.0 htmlwidgets_1.5.3
## [75] bit_4.0.4           tidyselect_1.1.1
## [77] magrittr_2.0.1      R6_2.5.0
## [79] gplots_3.1.1        generics_0.1.0
## [81] Hmisc_4.5-0         DelayedArray_0.19.0
## [83] DBI_1.1.1           pillar_1.6.1
## [85] foreign_0.8-81      survival_3.2-11
## [87] KEGGREST_1.33.0     RCurl_1.98-1.3
## [89] nnet_7.3-16         tibble_3.1.2
## [91] crayon_1.4.1        KernSmooth_2.23-20
## [93] utf8_1.2.1          BiocFileCache_2.1.0
## [95] jpeg_0.1-8.1        progress_1.2.2
## [97] grid_4.1.0          data.table_1.14.0
## [99] marray_1.71.0       blob_1.2.1
## [101] digest_0.6.27       munsell_0.5.0
## [103] Gviz_1.37.0

```