

msgbsR: an R package to analyse methylation sensitive genotyping by sequencing (MS-GBS) data

Benjamin Mayne

May 1, 2024

Contents

1	Introduction	2
2	Reading data into R	2
3	Confirmation of correct cut sites	3
4	Visualization of read counts	4
5	Differential methylation analysis	6
6	Visualization of cut site locations	6
7	Session Information	8
8	References	9

1 Introduction

Current data analysis tools do not fulfil all experimental designs. For example, GBS experiments using methylation sensitive restriction enzymes (REs), which is also known as methylation sensitive genotyping by sequencing (MS-GBS), is an effective method to identify differentially methylated sites that may not be accessible in other technologies such as microarrays and methyl capture sequencing. However, current data analysis tools do not satisfy the requirements for these types of experimental designs.

Here we present msgbsR, an R package for data analysis of MS-GBS experiments. Read counts and cut sites from a MS-GBS experiment can be read directly into the R environment from a sorted and indexed BAM file(s).

2 Reading data into R

The analysis with the msgbsR pipeline begins with a directory which contains sorted and indexed BAM file(s). msgbsR contains an example data set containing 6 samples from a MS-GBS experiment using the restriction enzyme MspI. In this example the 6 samples are from the prostate of a rat and have been truncated for chromosome 20. 3 of the samples were fed a control diet and the other 3 were fed an experimental high fat diet.

To read in the data directly into the R environment can be done using the `rawCounts()` function, which requires the directory path to where the sorted and indexed files are located and the desired number of threads to be run (Default = 1).

```
> library(msgbsR)
> library(GenomicRanges)
> library(SummarizedExperiment)
> my_path <- system.file("extdata", package = "msgbsR")
> se <- rawCounts(bamFilepath = my_path)
> dim(assay(se))
```

```
[1] 16047      6
```

The result is an `RangedSummarizedExperiment` object containing the read counts. The columns are samples and the rows contain the location of each unique cut sites. Each cut site has been given a unique ID (chr:position:position:strand). The cut site IDs can be turned into a `GRanges` object. Information regarding the samples such as treatment or other groups can be added into the return object as shown below

```
> colData(se) <- DataFrame(Group = c(rep("Control", 3), rep("Experimental", 3)),
+                             row.names = colnames(assay(se)))
```

3 Confirmation of correct cut sites

After the data has been generated into the R environment, the next step is to confirm that the cut sites were the correctly generated sites. In this example, the methylated sensitive restriction enzyme that has been used is MspI which recognizes a 4bp sequence (C/CGG). MspI cuts between the two cytosines when the outside cytosine is methylated.

The first step is to extract the location of the cut sites from `se` and adjust the cut sites such that the region will cover the recognition sequence of MspI. It is important to note that in this example the user must adjust the region over the cut sites specifically for each strand. In other words although the enzyme cuts at C/CGG on the minus strand this would appear as CCG/G. The code below shows how to adjust the postioning of the cut sites to cover the recognition site on each strand.

```
> cutSites <- rowRanges(se)
> # # Adjust the cut sites to overlap recognition site on each strand
> start(cutSites) <- ifelse(test = strand(cutSites) == '+',
+                           yes = start(cutSites) - 1, no = start(cutSites) - 2)
> end(cutSites) <- ifelse(test = strand(cutSites) == '+',
+                          yes = end(cutSites) + 2, no = end(cutSites) + 1)
```

The object `cutSites` is a `GRanges` object that contains the start and end position of the MspI sequence length around the cut sites. These cut sites can now be checked if the sequence matches the MspI sequence.

`msgbsR` offer two approaches to checking the cut sites. The first approach is to use a `BSgenome` which can be obtained from Bioconductor. In this example, `BSgenome.Rnorvegicus.UCSC.rn6` will be used.

```
> library(BSgenome.Rnorvegicus.UCSC.rn6)
> correctCuts <- checkCuts(cutSites = cutSites, genome = "rn6", seq = "CCGG")
```

If a `BSgenome` is unavailable for a species of interest, another option to checking the cut sites is to use a `fasta` file which can be used through the `checkCuts()` function.

The `correctCuts` data object is in the format of a `GRanges` object and contains the correct sites that contained the recognition sequence. These sites can be kept within `se` by using the `subsetByOverlaps` function.

The incorrect MspI cut sites can be filtered out of `datCounts`:

```
> se <- subsetByOverlaps(se, correctCuts)
> dim(assay(se))
```

```
[1] 13983    6
```

`se` now contains the correct cut sites and can now be used in downstream analyses.

4 Visualization of read counts

Before any further downstream analyses with the data, the user may want to filter out samples that did not generate a sufficient number of read counts or cut sites. The `msgbsR` package contains a function which plots the total number of read counts against the total number of cut sites produced per sample. The user can also use the function to visualize if different categories or groups produced varying amount of cut sites or total amount of reads.

To visualize the total number of read counts against the total number of cut sites produced per sample:

```
> plotCounts(se = se, category = "Group")
```

This function generates a plot (Figure 1) where the x axis and y axis represents the total number of reads and the total number of cut sites produced for each sample respectively.

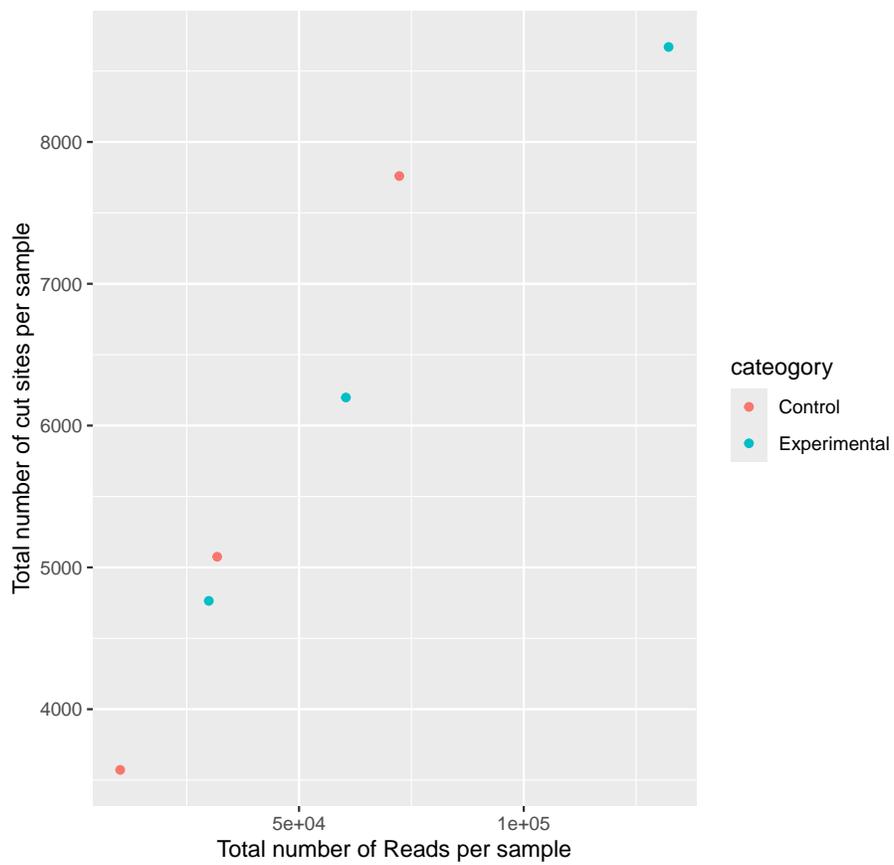


Figure 1: The distribution of the total number of reads and cut sites produced by each sample.

5 Differential methylation analysis

msgbsR utilizes edgeR in order to determine which cut sites are differentially methylated between groups. Since MS-GBS experiments can have multiple groups or conditions msgbsR offers a wrapper function of edgeR (Zhou et al., 2014) tools to automate differential methylation analyses.

To determine which cut sites are differentially methylated between groups:

```
> top <- diffMeth(se = se, category = "Group",
+               condition1 = "Control", condition2 = "Experimental",
+               cpmThreshold = 1, thresholdSamples = 1)
```

The top object now contains a data frame of the cut sites that had a CPM > 1 in at least 1 sample and which cut sites are differentially methylated between the two groups.

6 Visualization of cut site locations

The msgbsR package contains a function to allow visualization of the location of the cut sites. Given the lengths of the chromosomes the cut sites can be visualized in a circos plot (Figure 2).

Firstly, define the length of the chromosome.

```
> ratChr <- seqlengths(BSgenome.Rnorvegicus.UCSC.rn6)["chr20"]
```

Extract the differentially methylated cut sites.

```
> my_cuts <- GRanges(top$site[which(top$FDR < 0.05)])
```

To generate a circos plot:

```
> plotCircos(cutSites = my_cuts, seqlengths = ratChr,
+           cutSite.colour = "red", seqlengths.colour = "blue")
```

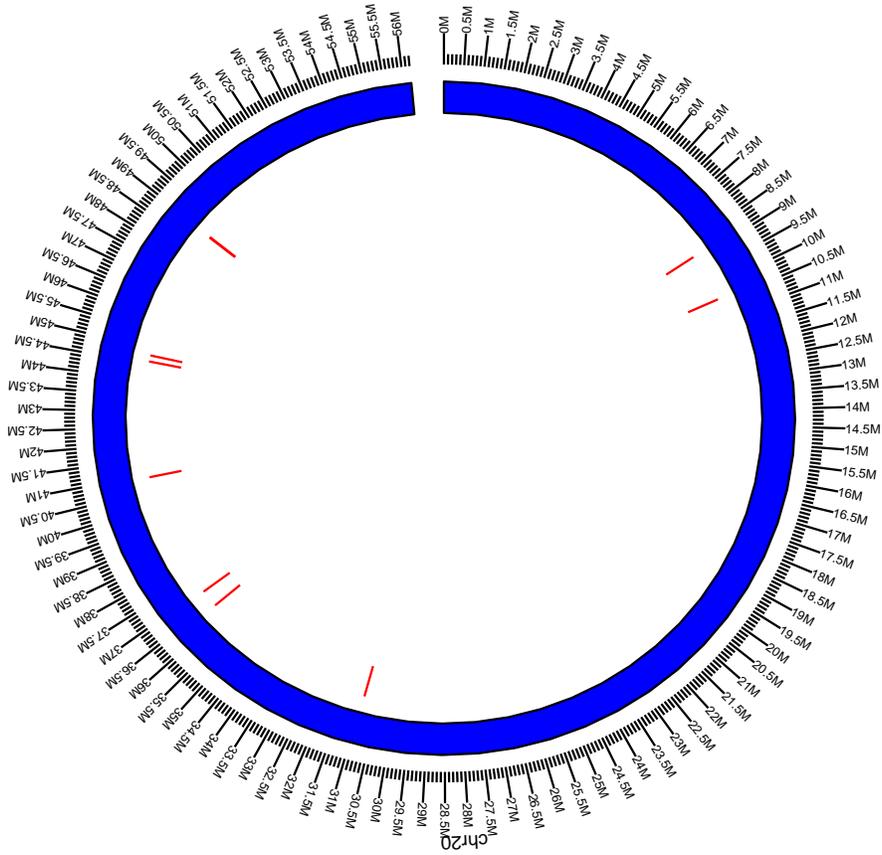


Figure 2: A circos plot of chromosome 20 representing cut sites defined by the user.

7 Session Information

This analysis was conducted on:

```
> sessionInfo()
```

```
R version 4.4.0 beta (2024-04-15 r86425)
```

```
Platform: x86_64-pc-linux-gnu
```

```
Running under: Ubuntu 22.04.4 LTS
```

```
Matrix products: default
```

```
BLAS: /home/biocbuild/bbs-3.19-bioc/R/lib/libRblas.so
```

```
LAPACK: /usr/lib/x86_64-linux-gnu/lapack/liblapack.so.3.10.0
```

```
locale:
```

```
[1] LC_CTYPE=en_US.UTF-8      LC_NUMERIC=C
[3] LC_TIME=en_US.UTF-8      LC_COLLATE=en_US.UTF-8
[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
```

```
time zone: America/New_York
```

```
tzcode source: system (glibc)
```

```
attached base packages:
```

```
[1] stats4      stats      graphics  grDevices  utils      datasets  methods
[8] base
```

```
other attached packages:
```

```
[1] BSgenome.Rnorvegicus.UCSC.rn6_1.4.1 BSgenome_1.72.0
[3] rtracklayer_1.64.0                   BiocIO_1.14.0
[5] Biostrings_2.72.0                     XVector_0.44.0
[7] SummarizedExperiment_1.34.0          Biobase_2.64.0
[9] MatrixGenerics_1.16.0                 matrixStats_1.3.0
[11] msgbsR_1.28.0                          GenomicRanges_1.56.0
[13] GenomeInfoDb_1.40.0                   IRanges_2.38.0
[15] S4Vectors_0.42.0                       BiocGenerics_0.50.0
```

```
loaded via a namespace (and not attached):
```

```
[1] splines_4.4.0                bitops_1.0-7                filelock_1.0.3
[4] tibble_3.2.1                  R.oo_1.26.0                 graph_1.82.0
[7] XML_3.99-0.16.1              rpart_4.1.23                lifecycle_1.0.4
[10] httr2_1.0.1                   pwalgn_1.0.0                edgeR_4.2.0
[13] lattice_0.22-6                ensemblDb_2.28.0            OrganismDbi_1.46.0
[16] backports_1.4.1              magrittr_2.0.3              limma_3.60.0
[19] Hmisc_5.1-2                   rmarkdown_2.26              yaml_2.3.8
```

[22] ggbio_1.52.0	DBI_1.2.2	RColorBrewer_1.1-3
[25] abind_1.4-5	ShortRead_1.62.0	zlibbioc_1.50.0
[28] purrr_1.0.2	R.utils_2.12.3	AnnotationFilter_1.28.0
[31] biovizBase_1.52.0	RCurl_1.98-1.14	nnet_7.3-19
[34] VariantAnnotation_1.50.0	rappdirs_0.3.3	GenomeInfoDbData_1.2.12
[37] codetools_0.2-20	DelayedArray_0.30.0	xml2_1.3.6
[40] tidyselect_1.2.1	UCSC.utils_1.0.0	farver_2.1.1
[43] BiocFileCache_2.12.0	base64enc_0.1-3	GenomicAlignments_1.40.0
[46] jsonlite_1.8.8	Formula_1.2-5	tools_4.4.0
[49] progress_1.2.3	Rcpp_1.0.12	glue_1.7.0
[52] gridExtra_2.3	SparseArray_1.4.0	xfun_0.43
[55] dplyr_1.1.4	genomeIntervals_1.60.0	withr_3.0.0
[58] BiocManager_1.30.22	fastmap_1.1.1	GGally_2.2.1
[61] latticeExtra_0.6-30	fansi_1.0.6	digest_0.6.35
[64] R6_2.5.1	colorspace_2.1-0	jpeg_0.1-10
[67] dichromat_2.0-0.1	biomaRt_2.60.0	RSQLite_2.3.6
[70] LSD_4.1-0	R.methodsS3_1.8.2	utf8_1.2.4
[73] tidyr_1.3.1	generics_0.1.3	intervals_0.15.4
[76] data.table_1.15.4	prettyunits_1.2.0	httr_1.4.7
[79] htmlwidgets_1.6.4	S4Arrays_1.4.0	ggstats_0.6.0
[82] pkgconfig_2.0.3	gtable_0.3.5	blob_1.2.4
[85] hwriter_1.3.2.1	htmltools_0.5.8.1	RBGL_1.80.0
[88] ProtGenerics_1.36.0	scales_1.3.0	png_0.1-8
[91] knitr_1.46	rstudioapi_0.16.0	reshape2_1.4.4
[94] rjson_0.2.21	checkmate_2.3.1	curl_5.2.1
[97] cachem_1.0.8	stringr_1.5.1	parallel_4.4.0
[100] foreign_0.8-86	AnnotationDbi_1.66.0	restfulr_0.0.15
[103] pillar_1.9.0	grid_4.4.0	vctrs_0.6.5
[106] easyRNASeq_2.39.0	dbplyr_2.5.0	cluster_2.1.6
[109] htmlTable_2.4.2	evaluate_0.23	GenomicFeatures_1.56.0
[112] cli_3.6.2	locfit_1.5-9.9	compiler_4.4.0
[115] Rsamtools_2.20.0	rlang_1.1.3	crayon_1.5.2
[118] labeling_0.4.3	interp_1.1-6	plyr_1.8.9
[121] stringi_1.8.3	deldir_2.0-4	BiocParallel_1.38.0
[124] txdbmaker_1.0.0	munsell_0.5.1	lazyeval_0.2.2
[127] Matrix_1.7-0	hms_1.1.3	bit64_4.0.5
[130] ggplot2_3.5.1	KEGGREST_1.44.0	statmod_1.5.0
[133] memoise_2.0.1	bit_4.0.5	

8 References

Zhou X, Lindsay H, Robinson MD (2014). Robustly detecting differential expression in RNA sequencing data using observation weights. *Nucleic Acids Research*, 42(11), e91.