

methyISig: A package for whole genome DNA methylation analysis

Yongseok Park, Maria E. Figueroa, Laura S. Rozek, and Maureen A. Sartor
email: yongpark@pitt.edu

June 17, 2014

Contents

1	Introduction	2
2	Installation	2
3	Basic	2
3.1	Loading methyISig package	2
3.2	Methylation score files	2
3.3	Reading methylation score files	3
4	Differential methylation analysis	3
4.1	Usage of the methyISigCalc function	3
4.2	Site specific analysis	4
4.3	Tiled Data analysis	4
4.4	Variance from one group	4
4.5	Using local information	4
5	Multiple threads computation	5
6	Annotation	5
6.1	CpG island Annotation	5
6.2	RefGene Annotation	6
6.3	Transcription factor enrichment test	7
7	Data visualization	7
8	Other	10
8.1	‘methyISigData’ object	10
8.1.1	S4 data structure	10
8.1.2	Subsetting	10
8.1.3	Getting values	11
8.2	‘methyISigDiff’ object	11
8.2.1	S4 data structure	11
8.2.2	Subsetting	11
8.2.3	Getting values	12
8.2.4	How to subtract DMCs or DMRs	12
8.3	Summarizing data	12
8.4	Generating heatmap	13
9	References	14

1 Introduction

DNA methylation plays critical roles in gene regulation and cellular specification without altering DNA sequences. It is one of the best understood and most intensively studied epigenetic marks in mammalian cells. Treatment of DNA with sodium bisulfite deaminates unmethylated cytosines to uracil while methylated cytosines are resistant to this conversion thus allowing for the discrimination between methylated and unmethylated CpG sites. Sodium bisulfite pre-treatment of DNA coupled with next-generation sequencing has allowed DNA methylation to be studied quantitatively and genome-wide at single cytosine site resolution.

MethylSig is a method for testing for differential methylated cytosines (DMCs) or regions (DMRs) in whole-genome bisulfite sequencing (bis-seq) or reduced representation bisulfite sequencing (RRBS) experiments. **MethylSig** uses a beta binomial model to test for significant differences between groups of samples. Several options exist for either site-specific or sliding window tests, combining strands, and for variance estimation. It allows annotating the resulting regions to multiple genome features, and visualizing the results for chosen genomic regions.

This document provides a step by step guide for the **methylSig** package.

2 Installation

MethylSig is available as an R package at <http://sartorlab.ccmb.med.umich.edu/software>. After downloading the **methylSig** package, for example, **methylSig_0.1.0.tar.gz**, the **install.packages()** function can be used to install **methylSig** in R. The current version of the **methylSig** package depends on R (≥ 2.10).

```
>install.packages("methylSig_0.1.0.tar.gz", repos=NULL, type="source")
* installing *source* package 'methylSig' ...
** R
** data
** inst
** preparing package for lazy loading
** help
*** installing help indices
** building package indices
** testing if installed package can be loaded
* DONE (methylSig)
```

3 Basic

3.1 Loading methylSig package

methylSig package can be loaded using the **library()** function.

```
>library(methylSig)
Loading required package: boot
Loading required package: parallel
```

This will load **methylSig** package as well as other dependent packages.

3.2 Methylation score files

CpG methylation score files can be obtained using programs such as **bismark**. The typical CpG methylation score files are as follows:

	chrBase	chr	base	strand	coverage	freqC	freqT
	chr21.43008527	chr21	43008527	F	32	100.00	0.00
	chr21.43008531	chr21	43008531	F	32	96.88	3.12
	chr21.43008543	chr21	43008543	F	32	90.62	9.38
	chr21.43008674	chr21	43008674	R	27	100.00	0.00
	chr21.43008710	chr21	43008710	R	67	94.03	5.97
	chr21.43008720	chr21	43008720	R	67	92.54	7.46

The CpG methylation score file must contain at least seven columns. Among these, second to seventh column must be, in order, chromosome, base, strand, coverage, percentage of Cs and percentage of Ts. Column names are not important. Strand format is F/R or +/-, where F/+ represents forward and R/- represents reverse strand.

3.3 Reading methylation score files

MethylSig package provides the `methylSigReadData()` function to read CpG methylation score files and convert these files into a 'methylSigData' object for further analysis and annotation.

```
> fileList = c(system.file("extdata", "AML_1.txt", package = "methylSig"),
               system.file("extdata", "AML_2.txt", package = "methylSig"),
               system.file("extdata", "AML_3.txt", package = "methylSig"),
               system.file("extdata", "AML_4.txt", package = "methylSig"),
               system.file("extdata", "NBM_1.txt", package = "methylSig"),
               system.file("extdata", "NBM_2.txt", package = "methylSig"),
               system.file("extdata", "NBM_3.txt", package = "methylSig"),
               system.file("extdata", "NBM_4.txt", package = "methylSig"))
> sample.id = c("AML1", "AML2", "AML3", "AML4", "NBM1", "NBM2", "NBM3", "NBM4")
> treatment = c(1,1,1,1,0,0,0,0)
> ##### Read Data #####
> meth <- methylSigReadData(fileList, sample.ids = sample.id, assembly = "hg18",
                           treatment = treatment, context = "CpG", destranded=TRUE)

Reading file (1/8) -- /tmp/RtmpYpYUhr/Rinst11f22c6f2b6/methylSig/extdata/AML_1.txt
Reading file (2/8) -- /tmp/RtmpYpYUhr/Rinst11f22c6f2b6/methylSig/extdata/AML_2.txt
Reading file (3/8) -- /tmp/RtmpYpYUhr/Rinst11f22c6f2b6/methylSig/extdata/AML_3.txt
Reading file (4/8) -- /tmp/RtmpYpYUhr/Rinst11f22c6f2b6/methylSig/extdata/AML_4.txt
Reading file (5/8) -- /tmp/RtmpYpYUhr/Rinst11f22c6f2b6/methylSig/extdata/NBM_1.txt
Reading file (6/8) -- /tmp/RtmpYpYUhr/Rinst11f22c6f2b6/methylSig/extdata/NBM_2.txt
Reading file (7/8) -- /tmp/RtmpYpYUhr/Rinst11f22c6f2b6/methylSig/extdata/NBM_3.txt
Reading file (8/8) -- /tmp/RtmpYpYUhr/Rinst11f22c6f2b6/methylSig/extdata/NBM_4.txt
(1)(2)(3)(4)(5)(6)(7)(8)
```

It is possible for the user to filter out CpG sites based on the read coverage. CpG sites with very large read coverage may be due to PCR bias and hence including CpG sites with very high coverage may distort the statistics of data analysis. The `methylSigReadData()` function provides 'minCount' and 'maxCount' arguments for defining lower and upper limits for coverage. The default values are 10 and 500 respectively.

```
> meth <- methylSigReadData(fileList, sample.ids = sample.id, assembly = "hg18",
                           treatment = treatment, context = "CpG", destranded=TRUE,
                           minCount=10, maxCount=500, quiet=TRUE)
```

Here, the argument `quiet` controls the progress reports. If `quiet=TRUE`, no progress reports will be shown.

There are many arguments for the `methylSigReadData()` function. Among these 'fileList', 'sample.ids' and 'treatment' are required. Some options have default values, for example, 'destranded=TRUE', 'num.cores=1', and 'quiet=FALSE'. Other arguments such as 'assembly', 'context' and 'pipeline' are optional and for information purposes only. The data type of `treatment` is a numeric vector. Each number represents a group. Multiple groups can be stored in one `methylSigData` object.

The argument 'num.cores' is used for multi-thread reading. See section 5 for more details.

4 Differential methylation analysis

4.1 Usage of the methylSigCalc function

The main function of this package is the differential methylation analysis function `methylSigCalc()`. It calculates differential methylation statistics between two groups of samples. It uses a beta-binomial approach to calculate differential methylation statistics, accounting for coverage and variation among samples within each group.

Usage:

```
methylSigCalc(meth, groups = c(Treatment = 1, Control = 0), dispersion = "both",  
              local.disp = FALSE, winsize.disp = 200, min.per.group = c(3, 3),  
              weightFunc=methylSig_weightFunc, num.cores = 1)
```

4.2 Site specific analysis

The default is to do site specific analysis and to use both groups to estimate variances.

```
> myDiffSigboth = methylSigCalc(meth, groups=c(1,0), min.per.group=3)
```

Total number of bases: 3.26k

The differentially methylated cytosines (DMCs) can be defined based on qvalues, pvalues and the methylation rate difference between two tested groups.

```
> myDiffSigbothDMCs = myDiffSigboth[myDiffSigboth[, "qvalue"] <= 0.05  
                                   & abs(myDiffSigboth[, "meth.diff"])>=25, ]
```

4.3 Tiled Data analysis

MethylSig package also provides `methylSigTile()` function to tile data within continuous non-overlapping windows. The default window size is 25bp. After tiling data, the `methylSigCalc()` function can be used to calculate differential methylation statistics.

```
> ### Tiled analysis  
> methTile = methylSigTile(meth, win.size = 25)  
> myDiffSigbothTile = methylSigCalc(methTile, groups=c(1,0), min.per.group=3)
```

Total number of regions: 994

4.4 Variance from one group

Using the `dispersion` argument, it is possible to estimate variances from one group rather than from both groups. The following code calculates differential methylation statistics based on estimating variances from group 0 only.

```
> ##### Variance from sample treatment group "0" only  
> myDiffSignorm = methylSigCalc(meth, groups=c(1,0), dispersion=0, min.per.group=3)
```

Total number of bases: 3.26k

```
> myDiffSignormTile = methylSigCalc(methTile, groups=c(1,0),  
                                   dispersion=0, min.per.group=3)
```

Total number of regions: 994

4.5 Using local information

It is also possible to use information from nearby CpG sites to improve the variance and methylation level estimates. The default `winsize.disp` and `winsize.meth` are 200 bps. The `winsize.disp` argument only takes into effect when `local.disp` is set to 'TRUE'. Similarly `winsize.meth` argument only takes into effect when `local.meth` is set to 'TRUE'.

```
> ##### Variance from both groups and using local information for variance  
> myDiffSigBothLoc = methylSigCalc(meth, groups=c(1,0),  
                                   min.per.group=3, local.disp=TRUE, winsize.disp=200)
```

Total number of bases: 3.26k

```
> ##### Variance from sample treatment group "0" only and using local information for variance  
> myDiffSignormLoc = methylSigCalc(meth, groups=c(1,0), dispersion=0,  
                                   min.per.group=3, local.disp=TRUE, winsize.disp=200)
```

Total number of bases: 3.26k

```
> ##### Variance from both groups and using local information for methylation level
> myDiffSigBothMLoc = methylSigCalc(meth, groups=c(1,0),
    min.per.group=3, local.meth=TRUE, winsize.meth=200)
```

Total number of bases: 3.26k

```
> ##### Variance from both groups and using local information for methylation level and variance
> myDiffSigBothMDLoc = methylSigCalc(meth, groups=c(1,0),
    min.per.group=3, local.disp=TRUE, winsize.disp=200,
    local.meth=TRUE, winsize.meth=200)
```

Total number of bases: 3.26k

5 Multiple threads computation

MethylSig package provides multi-thread programming to substantially reduce data analysis time. In the functions `methylSigDataRead()` and `methylSigCalc()`, multi-core programming will be initiated using 'num.cores' argument. Note that this option depends on R package 'parallel' and hence is not available in the Windows platform. The following example code is using 5 cores.

```
> #### Read Data using 5 cores
> meth <- methylSigReadData(fileList, sample.ids = sample.id, assembly = "hg18",
    treatment = treatment, context = "CpG", destranded=TRUE,
    num.cores=5, quiet=TRUE)
> #### Differential methylation analysis using 5 cores
> myDiffSigboth = methylSigCalc(meth, groups=c(1,0), min.per.group=3, num.cores=5)
```

Total number of bases: 3.26k

6 Annotation

6.1 CpG island Annotation

There are two functions, `cpgAnnotation()` and `cpgAnnotationPlot()`, in the `methylSig` package for CpG island annotation. The CpG island information file can be download from websites such as the UCSC genome browser. The appropriate genome assembly (the same genome assembly of the provided data) should be used.

In a linux server, the user may use the following command to download the annotation file for hg19. Please use appropriate directories for hg18, mm9 or mm10.

```
wget ftp://hgdownload.cse.ucsc.edu/goldenPath/hg19/database/cpgIslandExt.txt.gz
gunzip *.gz
```

Here we use the CpG island annotation file provided in the `methylSig` package to annotate our example. Note that this is a reduced annotation file and is not appropriate for a full real data analysis.

```
> library("graphics")
> cpgInfo = getCpGInfo(system.file("annotation", "cpgi.hg18.bed.txt",
    package = "methylSig"))
> myDiffDMCs = myDiffSigboth[myDiffSigboth[, "qvalue"] < 0.05
    & abs(myDiffSigboth[, "meth.diff"]) > 25, ]
> cpgAnn = cpgAnnotation(cpgInfo, myDiffSigboth)
> cpgAnnDmc = cpgAnnotation(cpgInfo, myDiffDMCs)
> cpgAnnotationPlot(cpgAnn, main="ALL")
> cpgAnnotationPlot(cpgAnnDmc, main="DMCs")
```

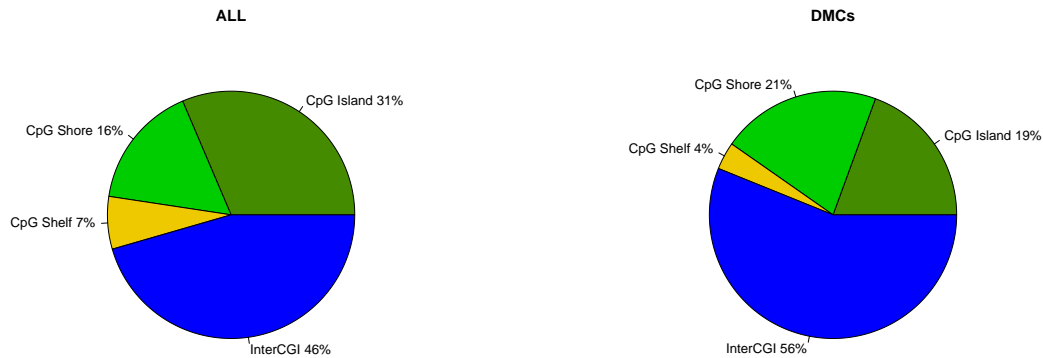


Figure 1: CpG annotation plot

6.2 RefGene Annotation

Again, there are two functions, `refGeneAnnotation()` and `refGeneAnnotationPlot()`, in `methyISig` package for annotation using RefGene models. The refGene information file can be download from websites such UCSC genome browser. The appropriate genome assembly (the same genome assembly of the provided data) should be used.

In a linux server, the user may use the following command to download the annotation file for hg19. Please use appropriate directories for hg18, mm9 or mm10.

```
wget ftp://hgdownload.cse.ucsc.edu/goldenPath/hg19/database/refGene.txt.gz
gunzip *.gz
```

We use refGene annotation file provided in the `methyISig` package to annotate in our example. Note that this is a reduced annotation file and is not appropriate for the a full real data analysis.

```
> refGeneInfo = getRefgeneInfo(system.file("annotation", "refGene.txt",
                                           package = "methyISig"))
> refGeneAnn = refGeneAnnotation(refGeneInfo, myDiffSigboth)
> refGeneAnnDmc = refGeneAnnotation(refGeneInfo, myDiffDMCs)
> refGeneAnnotationPlot(refGeneAnn, main="ALL",
                        priority=c("promoter", "cds", "noncoding", "5'utr", "3'utr"))
> refGeneAnnotationPlot(refGeneAnnDmc, main="DMC",
                        priority=c("promoter", "cds", "noncoding", "5'utr", "3'utr"))
```

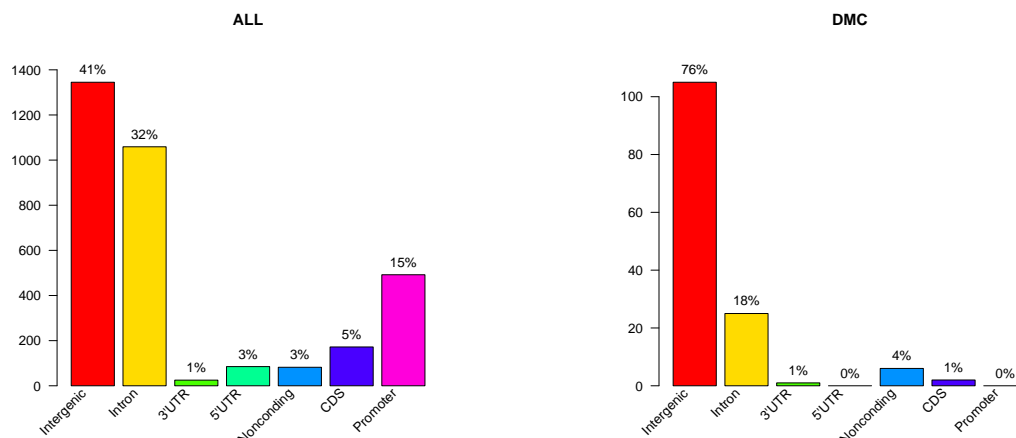


Figure 2: RefGene annotation plot

6.3 Transcription factor enrichment test

The functions `getTFBSInfo()` and `methyISig.tfbsEnrichTest()` are provided for reading the TFBS information file and implementing transcription factor enrichment test.

UCSC genome browser provides TFBS conserved track for hg18 and hg19. The following linux server shell command can be used to download these files:

```
wget ftp://hgdownload.cse.ucsc.edu/goldenPath/hg19/database/tfbsConsSites.txt.gz
wget ftp://hgdownload.cse.ucsc.edu/goldenPath/hg19/database/tfbsConsFactors.txt.gz
gunzip *.gz
```

Here, 'tfbsConsSites.txt' is tracking information and can be used directory in function `getTFBSInfo()`. The explanation of variable names is listed in file 'tfbsConsFactors.txt'.

Another TFBS track is from ENCODE for hg18, hg19 and mm9. However, the `methyISig` package cannot use this type of track directly. We provide ENCODE TFBS track files that are suitable for `methyISig` package at <http://sartorlab.ccmb.med.umich.edu/software>.

To identify which transcription factors (TFs) have significant level of hypermethylation or hypomethylation across their binding sites, which could indicate whether the TF is having a weaker or stronger regulatory effect, respectively, we first tile all reads from regions to which a particular TF is predicted to bind. We then apply our beta-binomial model to the data for each TF to identify TFs with hyper- or hypo-methylated binding sites.

To achieve this purpose, we provide fuction `methyISigTileTFBS()` to tile all data corresponding to the same TF.

```
> methTileTFs = methyISigTileTFBS(meth, tfbsInfo)
> myDiffTFs = methyISigCalc(methTileTFs, groups=c(1,0))
```

Total number of TFs: 126

7 Data visualization

MethyISig offers a unique two-tiered visualization of the methylation data depending on the zoom level. When the chromosome range is large (>1 million bp), the visualization function does not show individual sample data.

For narrow regions where at most 500 CpG sites have data reads, users can visualize sample-specific coverage levels and % methylation at each site, together with group averages, significance levels and a number of genomic annotations.

```

> tfbsInfo = getTFBSInfo(system.file("annotation", "tfbsUniform.txt",
                                     package = "methylSig"))
> DMCIndex = (myDiffSigboth[, "qvalue"] < 0.05
              & abs(myDiffSigboth[, "meth.diff"]) > 25)
> pvalue = methylSig.tfbsEnrichTest(myDiffSigboth, DMCIndex, tfbsInfo)

```

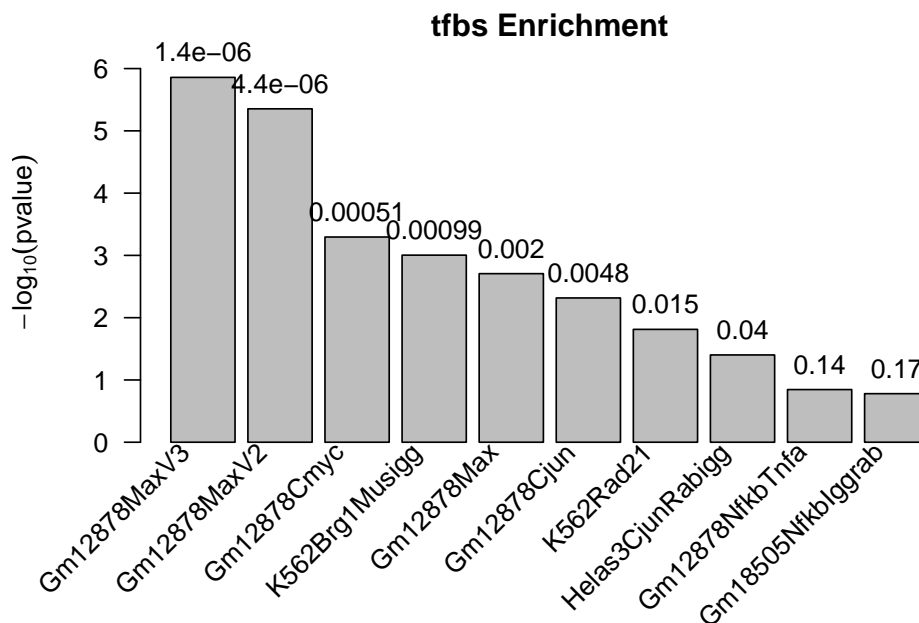


Figure 3: TFBS enrichment pvalue plot

```

> methylSigPlot(meth, "chr21", c(43000000, 44000000), groups=c(1,0),
               cpGInfo=cpGInfo, refGeneInfo=refGeneInfo,
               myDiff=myDiffSigboth, tfbsInfo=tfbsInfo, tfbsDense=F, sigQ=0.05)

```

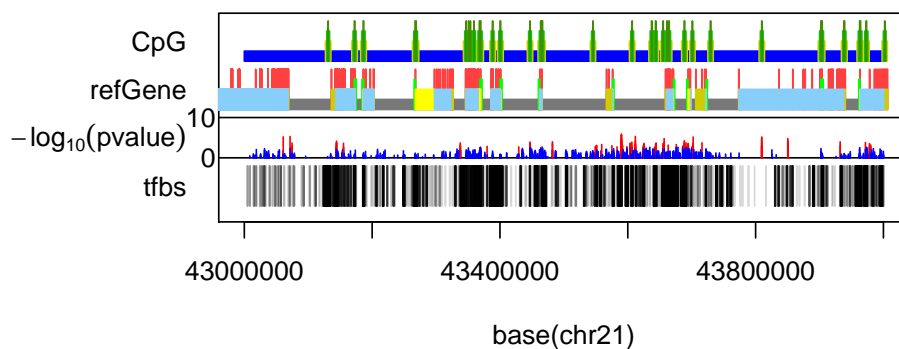


Figure 4: Data visualization in large range


```
> methylSigPlot(meth, "chr21", c(43800000, 43900000), groups=c(1,0),
  cpGInfo=cpGInfo, refGeneInfo=refGeneInfo,
  myDiff=myDiffSigboth, tfbsInfo=tfbsInfo, tfbsDense=F, sigQ=0.05)
```

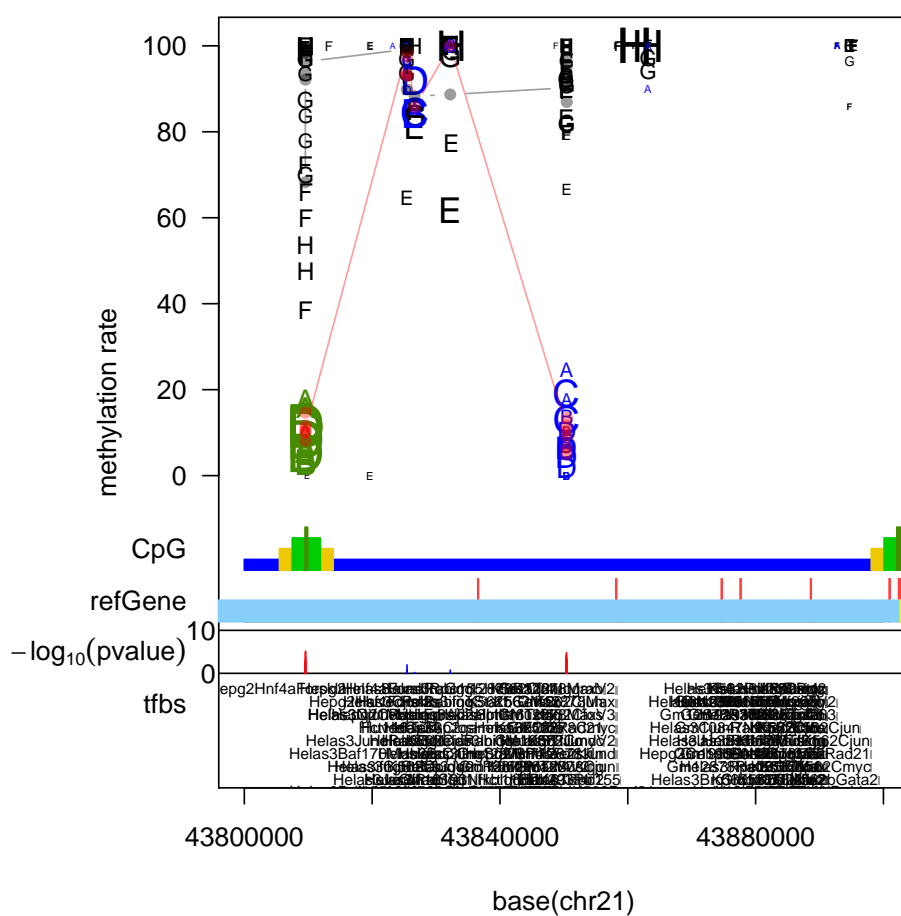


Figure 5: Data visualization within narrow range

8 Other

8.1 ‘methyISigData’ object

8.1.1 S4 data structure

the `methyISig` package uses S4 object. The contents of ‘methyISigData’ can be shown using the ‘`show()`’ function in R or just type the object itself.

```
> meth
```

```
methyISigData object with 7,571 rows
```

```
-----
      chr    start      end strand coverage1 numCs1 numTs1 ... coverage8 numCs8 numTs8
1 chr21 43000564 43000564    -      NA      NA      NA ...      NA      NA      NA
2 chr21 43000820 43000820    -      NA      NA      NA ...      NA      NA      NA
3 chr21 43008527 43008527    +      32      32      0 ...     124     122      2
4 chr21 43008531 43008531    +      32      31      1 ...     125     125      0
5 chr21 43008543 43008543    +      32      29      3 ...     125     117      8
6 chr21 43008665 43008665    -      NA      NA      NA ...      NA      NA      NA
7 chr21 43008673 43008673    -      27      27      0 ...      29      25      4
8 chr21 43008709 43008709    -      67      63      4 ...      31      26      5
9 chr21 43008719 43008719    -      67      62      5 ...      33      33      0
10 chr21 43014041 43014041    +      NA      NA      NA ...     207     202      5
-----
```

```
sample.ids: AML1 AML2 AML3 AML4 NBM1 NBM2 NBM3 NBM4
```

```
treatment: 1 1 1 1 0 0 0 0
```

```
destranded: TRUE
```

```
resolution: base
```

```
options: maxCount=500 & minCount=10 & assembly=hg18 & context=CpG
```

‘NA’ here means no data at this base location on the related sample.

8.1.2 Subsetting

Data can be subset using matrix style operations. Row represents base location and each column is a sample. Below is an example to obtain data for samples 1 to 4:

```
> meth1_4 = meth[,1:4]
```

This example returns the first 100 methylation reads in the data:

```
> methSub1_100 = meth[1:100,]
```

Two arguments can be used together. This example returns the first 100 methylation reads for samples 1 and 2.

```
> methSubData = meth[1:100,1:2]
```

```
> methSubData
```

```
methyISigData object with 60 rows
```

```
-----
      chr    start      end strand coverage1 numCs1 numTs1 coverage2 numCs2 numTs2
1 chr21 43008527 43008527    +      32      32      0      124     122      2
2 chr21 43008531 43008531    +      32      31      1      125     125      0
3 chr21 43008543 43008543    +      32      29      3      125     117      8
4 chr21 43008673 43008673    -      27      27      0      29      25      4
5 chr21 43008709 43008709    -      67      63      4      31      26      5
6 chr21 43008719 43008719    -      67      62      5      33      33      0
7 chr21 43014041 43014041    +      NA      NA      NA     207     202      5
8 chr21 43014044 43014044    +      NA      NA      NA     207     195     12
9 chr21 43014076 43014076    +      NA      NA      NA     202     188     14
-----
```

```
10 chr21 43016439 43016439      +      55      51      4      93      90      3
```

```
-----
sample.ids: AML1 AML2
treatment: 1 1
destranded: TRUE
resolution: base
options: maxCount=500 & minCount=10 & assembly=hg18 & context=CpG
```

8.1.3 Getting values

If the second argument is a string that matches one of the column names in the `methySigData` object, it gives the values of that column. Valid column names are “chr”, “start”, “end”, “strand”, “coverage1”, ..., “numCs1”, ..., and “numTs1”...

```
> coverage1 = meth[, "coverage1"]
> startTop200 = meth[1:200, "start"]
```

8.2 ‘methySigDiff’ object

8.2.1 S4 data structure

The contents of ‘methySigDiff’ are

```
> myDiffSigboth
```

methySigDiff object with 3,260 rows

```
-----
      chr    start    end strand    pvalue    qvalue  meth.diff logLikRatio
1  chr21 43008527 43008527      + 0.3219880 0.6664426 -2.4167475 1.27484057
2  chr21 43008531 43008531      + 0.8718018 1.0000000 0.4158129 0.02957797
3  chr21 43008543 43008543      + 0.5583045 0.8746145 -2.8071020 0.40674783
4  chr21 43014041 43014041      + 0.3670413 0.7056865 2.3585424 0.98274395
5  chr21 43014044 43014044      + 0.9166589 1.0000000 -0.6133413 0.01210852
6  chr21 43014076 43014076      + 0.1040236 0.4228389 12.1681339 4.39674527
7  chr21 43016439 43016439      + 0.2318347 0.5885776 -2.1713344 1.76894807
8  chr21 43016517 43016517      - 0.4054191 0.7342590 -2.0930416 0.86326517
9  chr21 43018213 43018213      + 0.8344303 1.0000000 -1.4186090 0.04765519
10 chr21 43018274 43018274      - 0.3425235 0.6784244 2.0009532 1.09884498
      theta df      mu1      mu0
1  3.689198e+01 4 96.49842 98.91517
2  1.172001e+01 4 98.45532 98.03950
3  3.131946e+01 4 90.97520 93.78231
4  2.829499e+01 5 97.82788 95.46934
5  1.269081e+01 5 90.45700 91.07034
6  1.904543e+01 4 93.23638 81.06825
7  9.151486e+01 6 96.96983 99.14116
8  1.000000e+06 4 93.10351 95.19656
9  7.094155e+00 6 91.07918 92.49778
10 7.115526e+01 5 98.16913 96.16818
-----
```

```
sample.ids: AML1 AML2 AML3 AML4 NBM1 NBM2 NBM3 NBM4
treatment: 1 1 1 1 0 0 0 0
destranded: TRUE
resolution: base
options: dispersion=both & local.disp=FALSE & local.meth=FALSE & min.per.group=c(3,3) & Total: 3260
```

8.2.2 Subsetting

This object can also subset by row to obtain results from part of CpG sites or regions. However, the qvalues will not be readjusted.

```
> myDiff100 = myDiffSigboth[1:100,]
```

8.2.3 Getting values

Similar to the 'methSigData' object, if the second argument is a string that is the same as one of the column names, it will return the results for that column. The valid variable names are "chr", "start", "end", "strand", "pvalue", "qvalue", "meth.diff", "logLikRatio", "theta", "df", "mu1", and "mu0". Here, for group methylation mean estimates "mu1" and "mu0", 1 and 0 come from the 'groups' argument in the methylSigCalc() function. So if one has run the methylSigCalc() function with 'groups=c(4,0)', then "mu4" and "mu0" will appear in the results.

```
> qvalues = myDiffSigboth[, "qvalue"]
```

8.2.4 How to subtract DMCs or DMRs

This 'methylSigDiff' object is very flexible to use by combining functions of subsetting and getting values. For example, the following code can obtain differentially methylated cytosines or regions defined as qvalue < 0.05 and difference of methylation rate > 25%.

```
> myDiffq05D25 = myDiffSigboth[myDiffSigboth[, "qvalue"] < 0.05  
                                & abs(myDiffSigboth[, "meth.diff"]) > 25,]
```

Here abs() is a R function to take an absolute value.

If you want to use pvalues instead of qvalues, then you can use

```
> myDiffp05D25 = myDiffSigboth[myDiffSigboth[, "pvalue"] < 0.05  
                                & abs(myDiffSigboth[, "meth.diff"]) > 25,]
```

8.3 Summarizing data

You can easily use other R functions to summarize or draw plots.

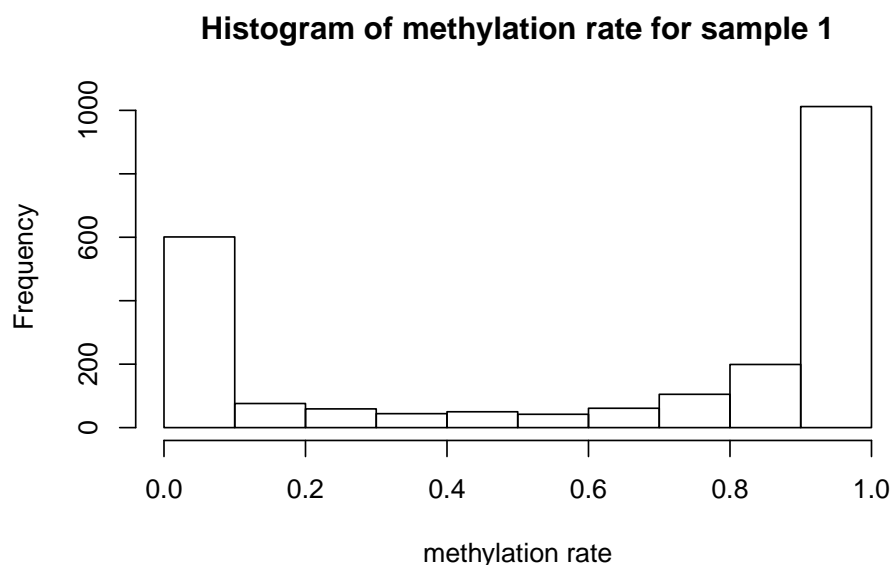
```
> methRaw = methylSigReadData(fileList, sample.ids = sample.id, assembly = "hg18",  
                             treatment = treatment, context = "CpG", minCount = 0,  
                             maxCount=Inf, destrand=F, quiet=T)  
> summary(methRaw[, "numCs1"]/methRaw[, "coverage1"])
```

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.	NA's
0.000	0.069	0.861	0.612	1.000	1.000	3606

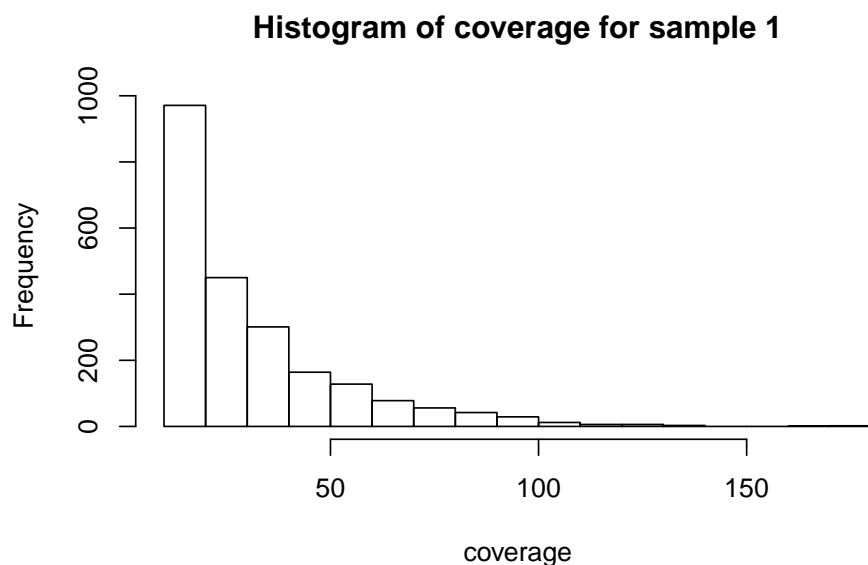
```
> summary(methRaw[, "coverage1"])
```

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.	NA's
10.00	14.00	23.00	30.68	39.00	177.00	3606

```
> hist(methRaw[, "numCs1"]/methRaw[, "coverage1"],  
       main="Histogram of methylation rate for sample 1",  
       xlab="methylation rate")
```



```
> hist(methRaw[, "coverage1"], main="Histogram of coverage for sample 1",
      xlab="coverage")
```



8.4 Generating heatmap

Here we provide an example to generate a correlation heatmap.

```
> library(gplots)
> x = meth[, "numCs"] / meth[, "coverage"]
> colnames(x) = meth@sample.ids
> rownames(x) = rep(NA, NROW(x))
> corrALL = cor(x, use="pairwise.complete.obs")
> heatmap.2(1-corrALL, na.rm=T, breaks=100,
  hclustfun = function(x) hclust(x, method="ward"),
  col="bluered", trace="none", symm=T, keysize=1, density.info="none")
```

Here is another example to generate a correlation heatmap based on differentially methylated cytosines.

```

> myDiffDMC = myDiffSigboth[myDiffSigboth[, "qvalue"] < 0.05
                        & abs(myDiffSigboth[, "meth.diff"]) >=25,]
> listInMeth = match(myDiffDMC@data.ids, meth@data.ids)
> y = x[listInMeth,]
> corrDMC = cor(y, use="pairwise.complete.obs")
> heatmap.2(1-corrDMC, na.rm=T, breaks=100,
            hclustfun = function(x) hclust(x,method="ward"),
            col="bluered",trace="none", symm=T, keysize=1,density.info="none")

```

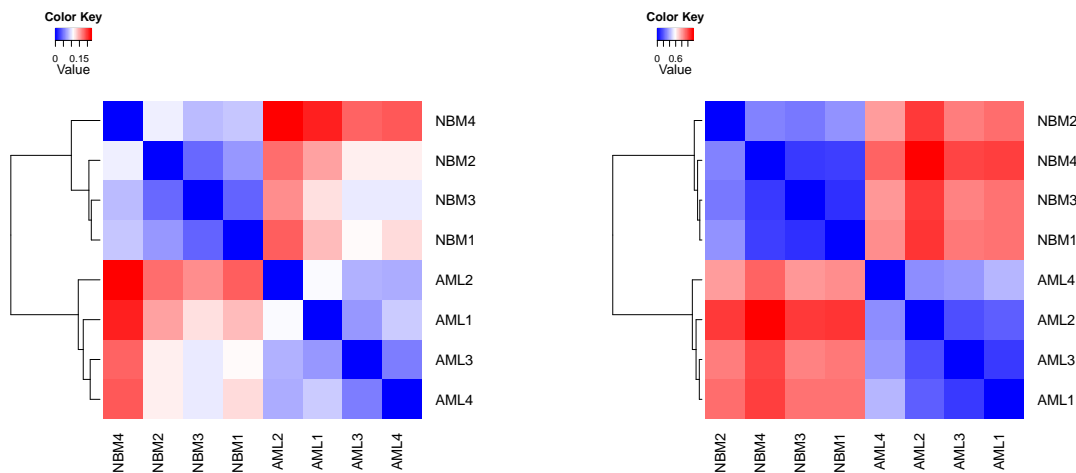


Figure 6: Heatmap based on methylation rate at all CpG sites (left) and at differentially methylated CpG sites (right)

9 References

Park, Y., Figueroa, M. E., Rozek, L. S., and sartor M.A., “methylSig: a whole genome DNA methylation analysis pipeline”. Submitted.