

The DMRcate package user's guide

Peters TJ

October 29, 2019

Summary

DMRcate extracts the most differentially methylated regions (DMRs) and variably methylated regions (VMRs) from both Whole Genome Bisulphite Sequencing (WGBS) and Illumina® Infinium BeadChip Array samples via kernel smoothing.

```
if (!require("BiocManager"))  
  install.packages("BiocManager")  
BiocManager::install("DMRcate")
```

Load DMRcate into the workspace:

```
library(DMRcate)
```

Illumina® Array Workflow

For this vignette, we will demonstrate DMRcate's array utility using data from ExperimentHub, namely Illumina HumanMethylationEPIC data from the data packages FlowSorted.Blood.EPIC. Specifically, we are interested in the methylation differences between CD4+ and CD8+ T cells.

```
library(ExperimentHub)  
eh <- ExperimentHub()  
FlowSorted.Blood.EPIC <- eh[["EH1136"]]  
tcell <- FlowSorted.Blood.EPIC[,colData(FlowSorted.Blood.EPIC)$CD4T==100 |  
                               colData(FlowSorted.Blood.EPIC)$CD8T==100]
```

For demonstration purposes, we will restrict this analysis to probes on chromosome 2:

```
locs <- IlluminaHumanMethylationEPICanno.ilm10b4.hg19::Locations  
chr2probes <- rownames(locs)[locs$chr=="chr2"]  
tcell <- subsetByLoci(tcell, includeLoci = chr2probes)  
  
## Loading required package: IlluminaHumanMethylationEPICmanifest
```

Firstly we have to filter out any probes where any sample has a failed position. Then we will normalise using `minfi::preprocessNoob`. After this, we extract the M -values from the `GenomicRatioSet`.

```
detP <- detectionP(tcell)
remove <- apply(detP, 1, function(x) any(x > 0.01))
tcell <- tcell[!remove,]
tcell <- preprocessNoob(tcell)
tcellms <- getM(tcell)
```

M-values (logit-transform of beta) are preferable to beta values for significance testing via `limma` because of increased sensitivity, but we will transform this to a beta matrix for visualisation purposes later on.

Some of the methylation measurements on the array may be confounded by proximity to SNPs, and cross-hybridisation to other areas of the genome [1, 2]. In particular, probes that are 0, 1, or 2 nucleotides from the methylcytosine of interest show a markedly different distribution to those farther away, in healthy tissue (Figure 1).

It is with this in mind that we filter out probes 2 nucleotides or closer to a SNP that have a minor allele frequency greater than 0.05, and the approximately 48,000 [1, 2] cross-reactive probes on either 450K and/or EPIC, so as to reduce confounding. Here we use a combination of *in silico* analyses from [1, 2]. About 4000 are removed from our M-matrix of approximately 65000:

```
nrow(tcellms)

## [1] 64714

tcellms.noSNPs <- rmSNPandCH(tcellms, dist=2, mafcut=0.05)
nrow(tcellms.noSNPs)

## [1] 60378
```

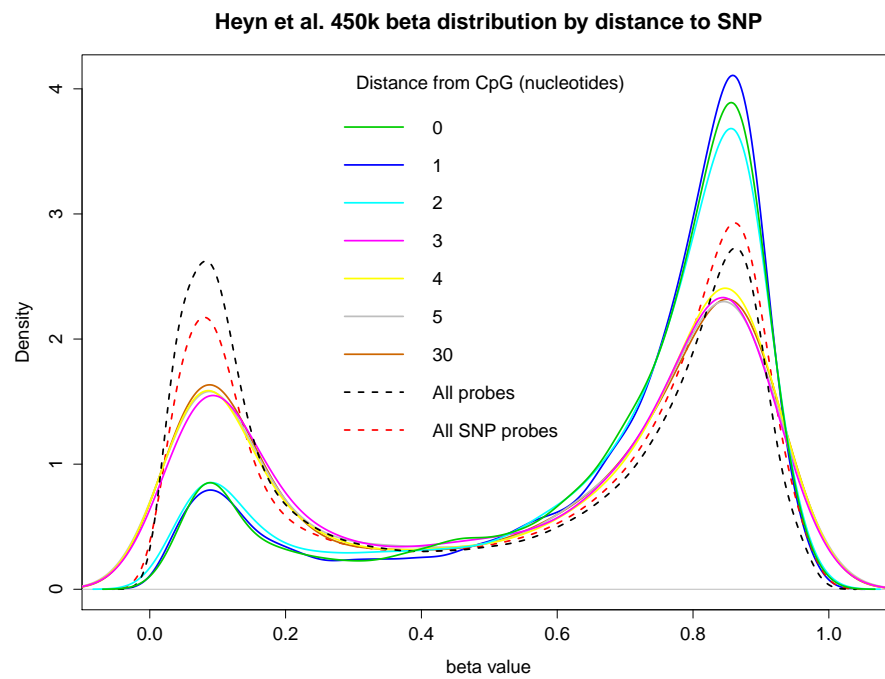
Here we have 6 CD8+ T cell assays, and 7 CD4+ T cell assays; we want to call DMRs between these groups. One of the CD4+ assays is a technical replicate, so we will average these two replicates like so:

```
tcell$Replicate

## [1] "" "" "" "" "" ""
## [7] "" "" "" "Th2535-1" "Th2535-1" ""
## [13] ""

tcell$Replicate[tcell$Replicate==""] <- tcell$Sample_Name[tcell$Replicate==""]
tcellms.noSNPs <- limma::avearrays(tcellms.noSNPs, tcell$Replicate)
tcell <- tcell[,!duplicated(tcell$Replicate)]
```

Figure 1: Beta distribution of 450K probes from publically available data from blood samples of healthy individuals [3] by their proximity to a SNP. “All SNP probes” refers to the 153 113 probes listed by Illumina® whose values may potentially be confounded by a SNP.



Next we want to annotate our matrix of M-values with relevant information. We also use the backbone of the `limma` pipeline for differential array analysis. We want to compare within patients across tissue samples, so we set up our variables for a standard `limma` pipeline, and set `coef=2` in `cpg.annotate` since this corresponds to the phenotype comparison in `design`.

`cpg.annotate()` takes either a data matrix with Illumina probe IDs, or an already prepared `GenomicRatioSet` from `minfi`.

```
type <- factor(tcell$CellType)
design <- model.matrix(~type)
myannotation <- cpg.annotate("array", tcellms.noSNPs, what="M", arraytype = "EPIC",
                             analysis.type="differential", design=design, coef=2)
```

```
myannotation

## CpGannotated object describing 60378 CpG sites, with independent
## CpG threshold indexed at fdr=0.05 and 2641 significant CpG sites.
```

Now we can find our most differentially methylated regions with `dmrcate()`.

For each chromosome, two smoothed estimates are computed: one weighted with per-CpG *t*-statistics and one not, for a null comparison. The two estimates are compared via a Satterthwaite approximation[4], and a significance test is calculated at all hg19 coordinates that an input probe maps to. After *fdr*-correction, regions are then agglomerated from groups of post-smoothed significant probes where the distance to the next consecutive probe is less than `lambda` nucleotides.

```
dmrcoutput <- dmrcate(myannotation, lambda=1000, C=2)

## Fitting chr2...
## Demarcating regions...
## Done!

dmrcoutput

## DMRResults object with 455 DMRs.
## Use extractRanges() to produce a GRanges object of these.
```

We can convert our DMR list to a `GRanges` object, which uses the `genome` argument to annotate overlapping gene loci.

```
results.ranges <- extractRanges(dmrcoutput, genome = "hg19")
results.ranges

## GRanges object with 455 ranges and 8 metadata columns:
##           seqnames           ranges strand |   no.cpgs
```

```

##          <Rle>          <IRanges> <Rle> | <integer>
##      [1]      chr2      87014979-87021117      * |      24
##      [2]      chr2 234294036-234297039      * |      14
##      [3]      chr2 112939119-112941244      * |       6
##      [4]      chr2   86991846-86992657      * |       3
##      [5]      chr2 197124443-197125372      * |       4
##      ...      ...      ...      ...      ...
## [451]      chr2 176986535-176987918      * |      13
## [452]      chr2 155553986-155555157      * |      15
## [453]      chr2   43454447-43455773      * |      21
## [454]      chr2   58654962-58655463      * |       7
## [455]      chr2 176964328-176964588      * |       7
##          min_smoothed_fdr          Stouffer          HMFDR
##          <numeric>          <numeric>          <numeric>
##      [1]          0 2.28490266103325e-58 3.09316274872618e-07
##      [2] 1.3550685101182e-98 1.3697237201002e-16 1.18753301105907e-05
##      [3] 2.04990646715572e-86 1.62119852947247e-17 2.24160948339874e-06
##      [4] 2.04557941055173e-186 1.03813997078552e-17 1.85819829548887e-07
##      [5] 5.3517222845669e-138 1.65186396027679e-17 8.92086803164279e-07
##      ...      ...      ...      ...
## [451] 7.32029615373132e-10 0.0213837630196166 0.24077103473153
## [452] 1.14577806363732e-11 0.0372538545821166 0.17800520785104
## [453] 1.12568888245577e-15 0.104822438189994 0.0454526910667972
## [454] 1.59068486346147e-09 0.360765088327784 0.093418658663056
## [455] 4.59616552500366e-09 0.105984279947368 0.228177727732663
##          Fisher          maxdiff          meandiff
##          <numeric>          <numeric>          <numeric>
##      [1] 3.11944098491378e-64 -0.73779911041198 -0.236569696749124
##      [2] 3.40435681423981e-19 -0.386692789729941 -0.137345011740568
##      [3] 2.85161930034251e-17 0.480622785916176 0.338850564829778
##      [4] 9.74218816837412e-17 -0.537150016399792 -0.398078510880661
##      [5] 1.48405871487718e-16 0.387780852388087 0.293760046300471
##      ...      ...      ...
## [451] 0.121730711855776 0.0778607668937892 0.0327760100181435
## [452] 0.126426770399768 0.0737463533211052 0.0240063380938431
## [453] 0.160161557332478 -0.0749531261442444 -0.0113009091171136
## [454] 0.224558284489568 0.0226519228285871 0.0109532578364987
## [455] 0.244599734093245 0.0691502531063555 0.0397062168865755
##
##
##      [1]          SNORA73, SNORA64, SNORA12, SNORA74, SNORA19, snR65, 5S_rRNA, SNORA4, SNORA73,
##      [2]
##      [3]          SNORA64, SNORA12, SNORA74, SNORA19, FBLN7, snR65, 5S_rRNA, SNORA4,
##      [4]          SNORA73, SNORA64, SNORA12, SNORA74, SNORA19, RMND5A, snR65, 5S_rRNA, SNORA4,
##      [5]

```

```
##      ...
##      [451]
##      [452]
##      [453]
##      [454] SNORA73, SNORA64, SNORD75, SNORA12, LINC01122, SNORA74, snR65, 5S_rRNA, SNORA4, S
##      [455]
##      -----
##      seqinfo: 1 sequence from an unspecified genome; no seqlengths
```

DMRs are ranked by Fisher's multiple comparison statistic, but **Stouffer** scores and the harmonic mean of the individual component FDRs (HMFDR) are also given in this object as alternative options for ranking DMR significance.

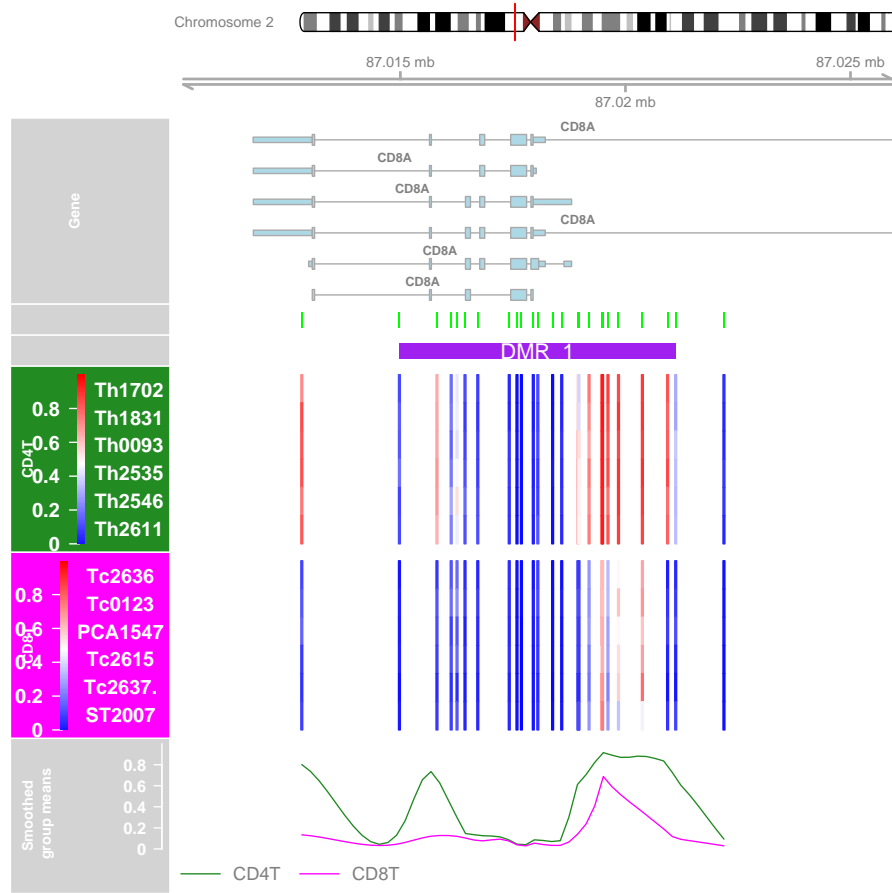
We can then pass this GRanges object to `DMR.plot()`, which uses the `Gviz` package as a backend for contextualising each DMR.

```
groups <- c(CD8T="magenta", CD4T="forestgreen")
cols <- groups[as.character(type)]
cols

##      CD4T      CD8T      CD8T      CD4T      CD4T
## "forestgreen" "magenta" "magenta" "forestgreen" "forestgreen"
##      CD8T      CD8T      CD8T      CD8T      CD4T
## "magenta" "magenta" "magenta" "magenta" "forestgreen"
##      CD4T      CD4T
## "forestgreen" "forestgreen"

tcellbetas <- ilogit2(tcellms.noSNPs)
colnames(tcellbetas) <- gsub("-1", "", tcell$Sample_Name)

DMR.plot(ranges=results.ranges, dmr=1, CpGs=tcellbetas, what="Beta",
         arraytype = "EPIC", phen.col=cols, genome="hg19")
```



Consonant with the expected biology, our top DMR shows the CD8+ T cells hypomethylated across parts of the CD8A locus. The two distinct hypomethylated sections have been merged because they are less than 1000 bp apart - specified by `lambda` in the call to `dmrcate()`. To call these as separate DMRs, make `lambda` smaller.

Bisulfite sequencing workflow

Bisulfite sequencing assays are fundamentally different to arrays, because methylation is represented as a pair of methylated and unmethylated reads per sample, instead of a single beta value. Although we could simply take the logit-proportion of methylated reads per CpG, this removes the effect of varying read depth across the genome. For example, a sampling depth of 30 methylated reads and 10 unmethylated reads is a much more precise estimate of the methylation level of a given CpG site than 3 methylated and 1 unmethylated. Hence, we take advantage of the fact that the overall effect can be expressed as an in-

teraction between the coefficient of interest and a two-level factor representing methylated and unmethylated reads [5].

The example shown here will be performed on a BSseq object containing bisulfite sequencing of regulatory T cells from various tissues as part of the `tissueTreg` package[6], imported using ExperimentHub. First, we will import the data:

```
bis_1072 <- eh[["EH1072"]]
bis_1072

## An object of type 'BSseq' with
## 21867550 methylation loci
## 15 samples
## has been smoothed with
## BSmooth (ns = 70, h = 1000, maxGap = 100000000)
## All assays are in-memory

colnames(bis_1072)

## [1] "Fat-Treg-R1"      "Fat-Treg-R2"      "Fat-Treg-R3"
## [4] "Liver-Treg-R1"    "Liver-Treg-R2"    "Liver-Treg-R3"
## [7] "Skin-Treg-R1"     "Skin-Treg-R2"     "Skin-Treg-R3"
## [10] "Lymph-N-Tcon-R1"  "Lymph-N-Tcon-R2"  "Lymph-N-Tcon-R3"
## [13] "Lymph-N-Treg-R1"  "Lymph-N-Treg-R2"  "Lymph-N-Treg-R3"
```

The data contains 15 samples: 3 (unmatched) replicates of mouse Tregs from fat, liver, skin and lymph node, plus a group of 3 CD4+ conventional lymph node T cells (Tcon). We will annotate the BSseq object to reflect this phenotypic information:

```
pData(bis_1072) <- data.frame(replicate=gsub(".*-", "", colnames(bis_1072)),
                             tissue=substr(colnames(bis_1072), 1,
                                             nchar(colnames(bis_1072))-3),
                             row.names=colnames(bis_1072))
colData(bis_1072)$tissue <- gsub("-", "_", colData(bis_1072)$tissue)
as.data.frame(colData(bis_1072))

##           replicate      tissue
## Fat-Treg-R1         R1  Fat_Treg
## Fat-Treg-R2         R2  Fat_Treg
## Fat-Treg-R3         R3  Fat_Treg
## Liver-Treg-R1        R1 Liver_Treg
## Liver-Treg-R2        R2 Liver_Treg
## Liver-Treg-R3        R3 Liver_Treg
## Skin-Treg-R1         R1  Skin_Treg
## Skin-Treg-R2         R2  Skin_Treg
```



```
## Skin-Treg-R3      R3      Skin_Treg
## Lymph-N-Tcon-R1   R1      Lymph_N_Tcon
## Lymph-N-Tcon-R2   R2      Lymph_N_Tcon
## Lymph-N-Tcon-R3   R3      Lymph_N_Tcon
## Lymph-N-Treg-R1   R1      Lymph_N_Treg
## Lymph-N-Treg-R2   R2      Lymph_N_Treg
## Lymph-N-Treg-R3   R3      Lymph_N_Treg
```

For standardisation purposes (and for `DMR.plot` to recognise the genome) we will change the chromosome naming convention to UCSC:

```
bis_1072 <- renameSeqlevels(bis_1072, mapSeqlevels(seqlevels(bis_1072), "UCSC"))
```

For demonstration purposes, we will retain CpGs on chromosome 19 only:

```
bis_1072 <- bis_1072[seqnames(bis_1072)=="chr19",]
bis_1072

## An object of type 'BSseq' with
##   558056 methylation loci
##   15 samples
## has been smoothed with
##   BSmooth (ns = 70, h = 1000, maxGap = 100000000)
## All assays are in-memory
```

Now we can prepare the model to be fit for `sequencing.annotate()`. The arguments are equivalent to `cpG.annotate()` but for a couple of exceptions:

- There is an extra argument `all.cov` giving an option whether to retain only CpGs where *all* samples have non-zero coverage, or whether to retain CpGs with only partial sample representation.
- The design matrix should be constructed to reflect the 2-factor structure of methylated and unmethylated reads. Fortunately, `edgeR::modelMatrixMeth()` can take a regular design matrix and transform it into the appropriate structure ready for model fitting.

```
tissue <- factor(pData(bis_1072)$tissue)
tissue <- relevel(tissue, "Liver_Treg")

#Regular matrix design
design <- model.matrix(~tissue)
colnames(design) <- gsub("tissue", "", colnames(design))
colnames(design)[1] <- "Intercept"
rownames(design) <- colnames(bis_1072)
design
```

```
##               Intercept Fat_Treg Lymph_N_Tcon Lymph_N_Treg Skin_Treg
## Fat-Treg-R1           1         1           0           0           0
## Fat-Treg-R2           1         1           0           0           0
## Fat-Treg-R3           1         1           0           0           0
## Liver-Treg-R1          1         0           0           0           0
## Liver-Treg-R2          1         0           0           0           0
## Liver-Treg-R3          1         0           0           0           0
## Skin-Treg-R1           1         0           0           0           1
## Skin-Treg-R2           1         0           0           0           1
## Skin-Treg-R3           1         0           0           0           1
## Lymph-N-Tcon-R1        1         0           1           0           0
## Lymph-N-Tcon-R2        1         0           1           0           0
## Lymph-N-Tcon-R3        1         0           1           0           0
## Lymph-N-Treg-R1        1         0           0           1           0
## Lymph-N-Treg-R2        1         0           0           1           0
## Lymph-N-Treg-R3        1         0           0           1           0
## attr("assign")
## [1] 0 1 1 1 1
## attr("contrasts")
## attr("contrasts")$tissue
## [1] "contr.treatment"

#Methylation matrix design
methdesign <- edgeR::modelMatrixMeth(design)
methdesign

##      Sample1 Sample2 Sample3 Sample4 Sample5 Sample6 Sample7 Sample8 Sample9
## 1           1         0         0         0         0         0         0         0
## 2           1         0         0         0         0         0         0         0
## 3           0         1         0         0         0         0         0         0
## 4           0         1         0         0         0         0         0         0
## 5           0         0         1         0         0         0         0         0
## 6           0         0         1         0         0         0         0         0
## 7           0         0         0         1         0         0         0         0
## 8           0         0         0         1         0         0         0         0
## 9           0         0         0         0         1         0         0         0
## 10          0         0         0         0         1         0         0         0
## 11          0         0         0         0         0         1         0         0
## 12          0         0         0         0         0         1         0         0
## 13          0         0         0         0         0         0         1         0
## 14          0         0         0         0         0         0         1         0
## 15          0         0         0         0         0         0         0         1
## 16          0         0         0         0         0         0         0         1
## 17          0         0         0         0         0         0         0         0
## 18          0         0         0         0         0         0         0         0
## 19          0         0         0         0         0         0         0         0
```

##	20	0	0	0	0	0	0	0	0	0
##	21	0	0	0	0	0	0	0	0	0
##	22	0	0	0	0	0	0	0	0	0
##	23	0	0	0	0	0	0	0	0	0
##	24	0	0	0	0	0	0	0	0	0
##	25	0	0	0	0	0	0	0	0	0
##	26	0	0	0	0	0	0	0	0	0
##	27	0	0	0	0	0	0	0	0	0
##	28	0	0	0	0	0	0	0	0	0
##	29	0	0	0	0	0	0	0	0	0
##	30	0	0	0	0	0	0	0	0	0
##		Sample10	Sample11	Sample12	Sample13	Sample14	Sample15	Intercept		
##	1	0	0	0	0	0	0	1		
##	2	0	0	0	0	0	0	0		
##	3	0	0	0	0	0	0	1		
##	4	0	0	0	0	0	0	0		
##	5	0	0	0	0	0	0	1		
##	6	0	0	0	0	0	0	0		
##	7	0	0	0	0	0	0	1		
##	8	0	0	0	0	0	0	0		
##	9	0	0	0	0	0	0	1		
##	10	0	0	0	0	0	0	0		
##	11	0	0	0	0	0	0	1		
##	12	0	0	0	0	0	0	0		
##	13	0	0	0	0	0	0	1		
##	14	0	0	0	0	0	0	0		
##	15	0	0	0	0	0	0	1		
##	16	0	0	0	0	0	0	0		
##	17	0	0	0	0	0	0	1		
##	18	0	0	0	0	0	0	0		
##	19	1	0	0	0	0	0	1		
##	20	1	0	0	0	0	0	0		
##	21	0	1	0	0	0	0	1		
##	22	0	1	0	0	0	0	0		
##	23	0	0	1	0	0	0	1		
##	24	0	0	1	0	0	0	0		
##	25	0	0	0	1	0	0	1		
##	26	0	0	0	1	0	0	0		
##	27	0	0	0	0	1	0	1		
##	28	0	0	0	0	1	0	0		
##	29	0	0	0	0	0	1	1		
##	30	0	0	0	0	0	1	0		
##		Fat_Treg	Lymph_N_Tcon	Lymph_N_Treg	Skin_Treg					
##	1	1	0	0	0					
##	2	0	0	0	0					

```
## 3      1      0      0      0
## 4      0      0      0      0
## 5      1      0      0      0
## 6      0      0      0      0
## 7      0      0      0      0
## 8      0      0      0      0
## 9      0      0      0      0
## 10     0      0      0      0
## 11     0      0      0      0
## 12     0      0      0      0
## 13     0      0      0      1
## 14     0      0      0      0
## 15     0      0      0      1
## 16     0      0      0      0
## 17     0      0      0      1
## 18     0      0      0      0
## 19     0      1      0      0
## 20     0      0      0      0
## 21     0      1      0      0
## 22     0      0      0      0
## 23     0      1      0      0
## 24     0      0      0      0
## 25     0      0      1      0
## 26     0      0      0      0
## 27     0      0      1      0
## 28     0      0      0      0
## 29     0      0      1      0
## 30     0      0      0      0
```

Just like for `cpg.annotate()`, we can specify a contrast matrix to find our comparisons of interest.

```
cont.mat <- limma::makeContrasts(treg_vs_tcon=Lymph_N_Treg-Lymph_N_Tcon,
                                fat_vs_ln=Fat_Treg-Lymph_N_Treg,
                                skin_vs_ln=Skin_Treg-Lymph_N_Treg,
                                fat_vs_skin=Fat_Treg-Skin_Treg,
                                levels=methdesign)

cont.mat
```

##	Contrasts	treg_vs_tcon	fat_vs_ln	skin_vs_ln	fat_vs_skin
## Levels					
## Sample1		0	0	0	0
## Sample2		0	0	0	0
## Sample3		0	0	0	0
## Sample4		0	0	0	0
## Sample5		0	0	0	0

##	Sample6	0	0	0	0
##	Sample7	0	0	0	0
##	Sample8	0	0	0	0
##	Sample9	0	0	0	0
##	Sample10	0	0	0	0
##	Sample11	0	0	0	0
##	Sample12	0	0	0	0
##	Sample13	0	0	0	0
##	Sample14	0	0	0	0
##	Sample15	0	0	0	0
##	Intercept	0	0	0	0
##	Fat_Treg	0	1	0	1
##	Lymph_N_Tcon	-1	0	0	0
##	Lymph_N_Treg	1	-1	-1	0
##	Skin_Treg	0	0	1	-1

Say we want to find DMRs between the regulatory and conventional T cells from the lymph node. First we would fit the model, where `sequencing.annotate()` transforms counts into log2CPMs (via `limma::voom()`) and uses `limma` under the hood to generate per-CpG *t*-statistics, indexing the FDR at 0.05:

```
seq_annot <- sequencing.annotate(bis_1072, methdesign, all.cov = TRUE,
                                contrasts = TRUE, cont.matrix = cont.mat,
                                coef = "treg_vs_tcon", fdr=0.05)

## Filtering out all CpGs where at least one sample has zero coverage...
## Processing BSseq object...
## Transforming counts...
## Fitting model...
## Your contrast returned 157 individually significant CpGs. We recommend
## the default setting of pcutoff in dmrcate().

seq_annot

## CpGannotated object describing 506908 CpG sites, with independent
## CpG threshold indexed at fdr=0.05 and 157 significant CpG sites.
```

And then, just like before, we can call DMRs with `dmrcate()`:

```
dmrcate.res <- dmrcate(seq_annot, C=2, min.cpgs = 5)

## Fitting chr19...
## Demarcating regions...
## Done!

dmrcate.res
```

```

## DMRResults object with 9 DMRs.
## Use extractRanges() to produce a GRanges object of these.

treg_vs_tcon.ranges <- extractRanges(dmrcate.res, genome="mm10")

## snapshotDate(): 2019-10-22
## see ?DMRcatedata and browseVignettes('DMRcatedata') for documentation
## loading from cache

treg_vs_tcon.ranges

## GRanges object with 9 ranges and 8 metadata columns:
##           seqnames           ranges strand |   no.cpgs   min_smoothed_fdr
##           <Rle>             <IRanges> <Rle> | <integer>         <numeric>
## [1]   chr19 29270611-29272005      * |      16 4.32351382071251e-94
## [2]   chr19 26683453-26684174      * |      12 1.77927194052734e-57
## [3]   chr19 32276886-32278089      * |      13 1.74619734491989e-56
## [4]   chr19 29374953-29375393      * |      12 1.48256678096532e-54
## [5]   chr19 36378257-36379597      * |      27 1.53747431922626e-76
## [6]   chr19 46653280-46654180      * |      19 3.94008431277526e-59
## [7]   chr19 57092365-57092646      * |      10 3.80467545599821e-36
## [8]   chr19 40808208-40809554      * |      26 3.43872692336671e-63
## [9]   chr19 41874401-41874895      * |      22 2.75828631295851e-39
##           Stouffer           HMFDR           Fisher
##           <numeric>         <numeric>         <numeric>
## [1]           1 0.0151786468767818 2.14645276178178e-08
## [2]           1 0.00787739676169919 0.000128162777894713
## [3]           1 0.0446758599711491 0.000150766828389467
## [4]           1 0.028226547829304 0.00241190715810346
## [5]           1 0.0482585210991692 0.00725026464290328
## [6]           1 0.0512002475534357 0.0452566259528858
## [7] 0.139493604737052 0.0711192891942912 0.0639020544512846
## [8]           1 0.180257052466182 0.305279374376655
## [9]           1 0.185853374084263 0.690216893652041
##           maxdiff           meandiff overlapping.genes
##           <numeric>         <numeric>         <character>
## [1] -6.40482000070317 -4.22352877428813      Jak2
## [2] -6.4032835829523 -3.53692271812483      Smarca2
## [3]  5.81469634649427  3.93201028394053      Sgms1
## [4] -6.10902321908482 -3.02082671194158 Cd274, AC119228.1
## [5] -6.09624814631612 -3.0355027665284      Pcgf5
## [6]  5.1838807268792  2.93151799486439      Wbp11
## [7] -4.67644961940502 -3.36472486679425      Ablim1
## [8] -4.83855429451068 -3.0749438164254      Cc2d2b
## [9]  4.57010787202641  2.56520038289291      Rrp12
## -----
## seqinfo: 1 sequence from an unspecified genome; no seqlengths

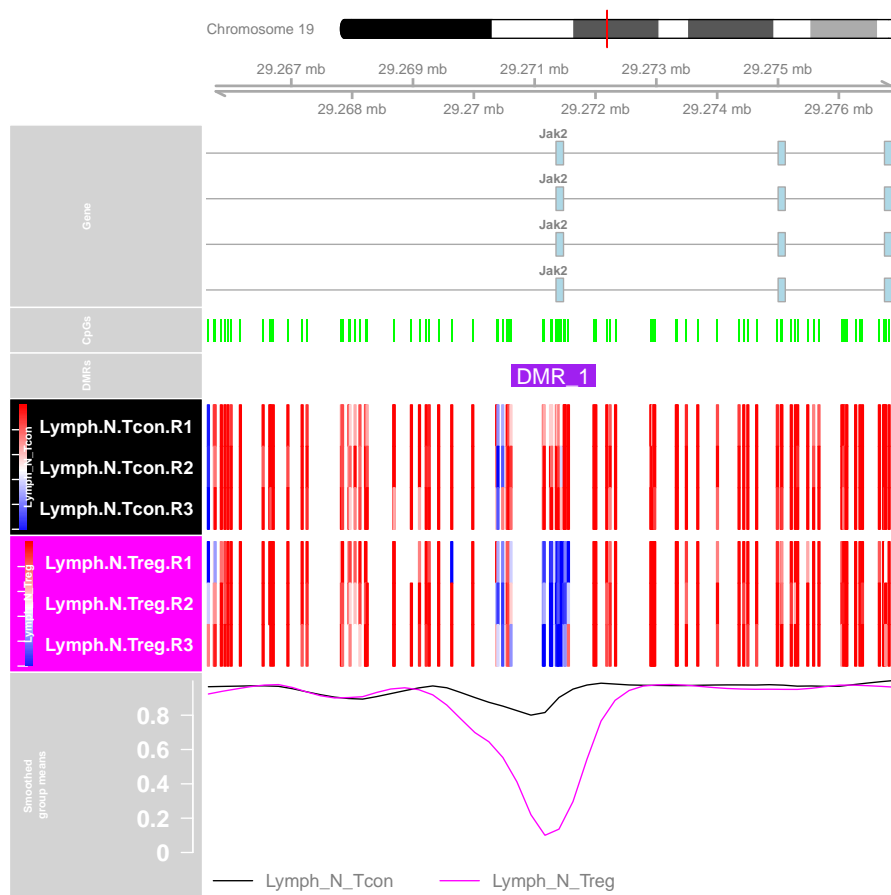
```

Looks like the top DMR is associated with the *Jak2* locus and hypomethylation in the Treg cells (since `meandiff < 0`). We can plot it like so:

```
cols <- as.character(plyr::mapvalues(tissue, unique(tissue),
                                     c("darkorange", "maroon", "blue",
                                       "black", "magenta")))

names(cols) <- tissue

DMR.plot(treg_vs_tcon.ranges, dmr = 1,
         CpGs=bis_1072[,tissue %in% c("Lymph_N_Tcon", "Lymph_N_Treg")],
         phen.col = cols[tissue %in% c("Lymph_N_Tcon", "Lymph_N_Treg")],
         genome="mm10")
```



Now, let's find DMRs between fat and skin Tregs.

```
seq_annot <- sequencing.annotate(bis_1072, methdesign, all.cov = TRUE,
                                contrasts = TRUE, cont.matrix = cont.mat,
                                coef = "fat_vs_skin", fdr=0.05)

## Filtering out all CpGs where at least one sample has zero coverage...
## Processing BSseq object...
## Transforming counts...
## Fitting model...
## Your contrast returned 5 individually significant CpGs; a small
but real effect. Consider increasing the 'fdr' parameter using changeFDR(),
but be warned there is an increased risk of Type I errors.
```

Because this comparison is a bit more subtle, there are very few significantly differential CpGs at this threshold. So we can use `changeFDR()` to relax the FDR to 0.25, taking into account that there is an increased risk of false positives.

```
seq_annot <- changeFDR(seq_annot, 0.25)

## Threshold is now set at FDR=0.25, resulting in 63 significantly differential CpGs.
```

```
dmrcate.res <- dmrcate(seq_annot, C=2, min.cpgs = 5)

## Fitting chr19...
## Demarcating regions...
## Done!

fat_vs_skin.ranges <- extractRanges(dmrcate.res, genome="mm10")

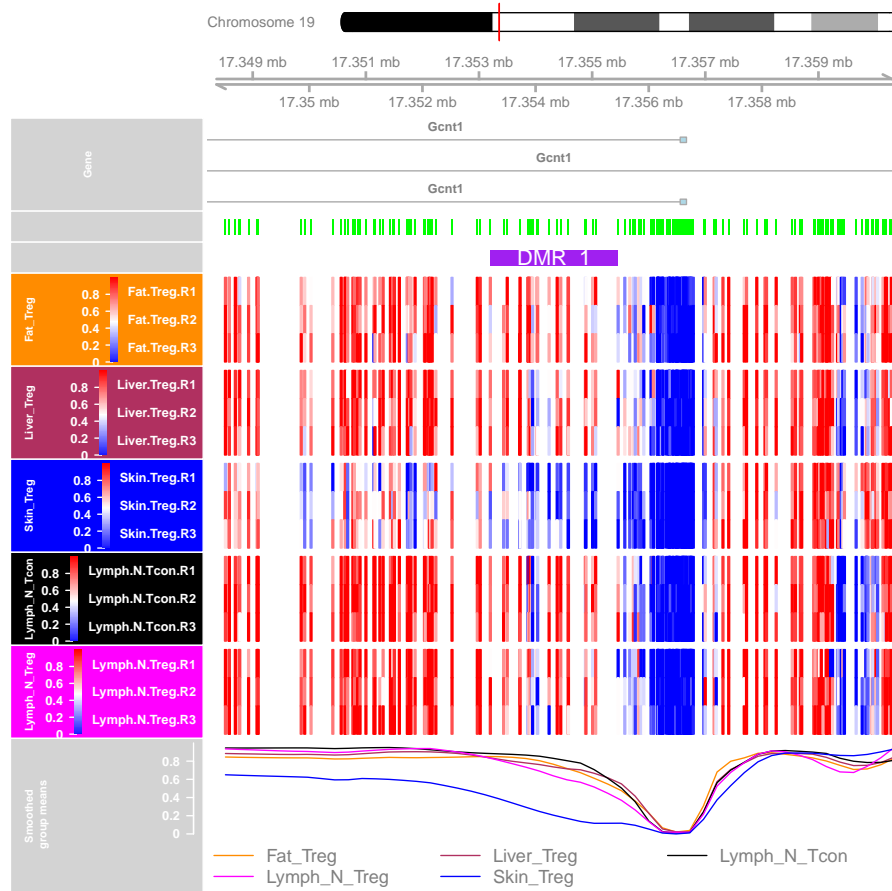
## snapshotDate(): 2019-10-22
## see ?DMRcatedata and browseVignettes('DMRcatedata') for documentation
## loading from cache
```

Now let's plot the top DMR with not only fat and skin, but with all samples:

```
cols

##      Fat_Treg      Fat_Treg      Fat_Treg      Liver_Treg      Liver_Treg
## "darkorange" "darkorange" "darkorange"      "maroon"      "maroon"
##      Liver_Treg      Skin_Treg      Skin_Treg      Skin_Treg      Lymph_N_Tcon
##      "maroon"      "blue"      "blue"      "blue"      "black"
## Lymph_N_Tcon Lymph_N_Tcon Lymph_N_Treg Lymph_N_Treg Lymph_N_Treg
##      "black"      "black"      "magenta"      "magenta"      "magenta"

DMR.plot(fat_vs_skin.ranges, dmr = 1, CpGs=bis_1072, phen.col = cols, genome="mm10")
```

Here we can see the methylation of skin cells over this section of *Gcnt1* is hypomethylated not only relative to fat, but to the other tissues as well.

As an alternative to `limma`, there is also the option of taking CpG-level differential statistics using `DSS::DMLtest()` or `DSS::DMLtest.multiFactor()`. There is no need to pass arguments such as `design`, `coef`, etc. to `sequencing.annotate()` in this case since we do this outside of the function. `fdr`, however, must be specified. For example:

```
library(DSS)
DMLfit <- DMLfit.multiFactor(bis_1072, design=data.frame(tissue=tissue), formula=~tissue)

## Fitting DML model for CpG site: 100000 , 200000 , 300000 , 400000 , 500000 ,

DSS_treg.vs.tcon <- DMLtest.multiFactor(DMLfit, Contrast=matrix(c(0, 0, -1, 1, 0)))
#Make sure to filter out all sites where the test statistic is NA
DSS_treg.vs.tcon <- DSS_treg.vs.tcon[!is.na(DSS_treg.vs.tcon$stat),]
```

```

seq_annot <- sequencing.annotate(obj=DSS_treg.vs.tcon, fdr=0.05)
seq_annot

## CpGannotated object describing 544489 CpG sites, with independent
## CpG threshold indexed at fdr=0.05 and 450 significant CpG sites.

dmrcate.res <- dmrcate(seq_annot, C=2, min.cpgs = 5)
DSS.treg_vs_tcon.ranges <- extractRanges(dmrcate.res, genome="mm10")

findOverlaps(treg_vs_tcon.ranges, DSS.treg_vs_tcon.ranges)

## Hits object with 9 hits and 0 metadata columns:
##      queryHits subjectHits
##      <integer>  <integer>
## [1]          1          1
## [2]          2          3
## [3]          3          5
## [4]          4          9
## [5]          5          2
## [6]          6         15
## [7]          7         26
## [8]          8         18
## [9]          9         24
## -----
## queryLength: 9 / subjectLength: 30

```

All of the 9 DMRs found using results from `limma` are also found using `DSS::DMLtest.multiFactor()`, with an extra 21 DMRs found by the latter at the same FDR. This suggests that `DMLtest.multiFactor()` is a little more permissive in calling differential methylation.

```

sessionInfo()

## R version 3.6.1 (2019-07-05)
## Platform: x86_64-pc-linux-gnu (64-bit)
## Running under: Ubuntu 18.04.3 LTS
##
## Matrix products: default
## BLAS:   /home/biocbuild/bbs-3.10-bioc/R/lib/libRblas.so
## LAPACK: /home/biocbuild/bbs-3.10-bioc/R/lib/libRlapack.so
##
## locale:
##  [1] LC_CTYPE=en_US.UTF-8      LC_NUMERIC=C
##  [3] LC_TIME=en_US.UTF-8      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] splines stats4 parallel stats graphics grDevices utils
## [8] datasets methods base
##
## other attached packages:
## [1] DSS_2.34.0
## [2] bsseq_1.22.0
## [3] tissueTreg_1.5.0
## [4] DMRcatedata_1.99.0
## [5] IlluminaHumanMethylationEPICmanifest_0.3.0
## [6] FlowSorted.Blood.EPIC_1.3.0
## [7] IlluminaHumanMethylationEPICanno.ilm10b4.hg19_0.6.0
## [8] nlme_3.1-141
## [9] quadprog_1.5-7
## [10] genefilter_1.68.0
## [11] ExperimentHub_1.12.0
## [12] AnnotationHub_2.18.0
## [13] BiocFileCache_1.10.0
## [14] dbplyr_1.4.2
## [15] DMRcate_2.0.0
## [16] minfi_1.32.0
## [17] bumphunter_1.28.0
## [18] locfit_1.5-9.1
## [19] iterators_1.0.12
## [20] foreach_1.4.7
## [21] Biostrings_2.54.0
## [22] XVector_0.26.0
## [23] SummarizedExperiment_1.16.0
## [24] DelayedArray_0.12.0
## [25] BiocParallel_1.20.0
## [26] matrixStats_0.55.0
## [27] Biobase_2.46.0
## [28] GenomicRanges_1.38.0
## [29] GenomeInfoDb_1.22.0
## [30] IRanges_2.20.0
## [31] S4Vectors_0.24.0
## [32] BiocGenerics_0.32.0
##
## loaded via a namespace (and not attached):
## [1] R.utils_2.9.0
## [2] tidyselect_0.2.5
## [3] RSQLite_2.1.2

```

```
## [4] AnnotationDbi_1.48.0
## [5] htmlwidgets_1.5.1
## [6] grid_3.6.1
## [7] munsell_0.5.0
## [8] codetools_0.2-16
## [9] preprocessCore_1.48.0
## [10] statmod_1.4.32
## [11] withr_2.1.2
## [12] colorspace_1.4-1
## [13] highr_0.8
## [14] knitr_1.25
## [15] rstudioapi_0.10
## [16] GenomeInfoDbData_1.2.2
## [17] bit64_0.9-7
## [18] rhdf5_2.30.0
## [19] vctrs_0.2.0
## [20] xfun_0.10
## [21] biovizBase_1.34.0
## [22] R6_2.4.0
## [23] illuminaio_0.28.0
## [24] AnnotationFilter_1.10.0
## [25] bitops_1.0-6
## [26] reshape_0.8.8
## [27] assertthat_0.2.1
## [28] promises_1.1.0
## [29] IlluminaHumanMethylation450kanno.ilmn12.hg19_0.6.0
## [30] scales_1.0.0
## [31] nnet_7.3-12
## [32] gtable_0.3.0
## [33] methylumi_2.32.0
## [34] ensemblDb_2.10.0
## [35] rlang_0.4.1
## [36] zeallot_0.1.0
## [37] rtracklayer_1.46.0
## [38] lazyeval_0.2.2
## [39] acepack_1.4.1
## [40] GEOquery_2.54.0
## [41] dichromat_2.0-0
## [42] checkmate_1.9.4
## [43] BiocManager_1.30.9
## [44] yaml_2.2.0
## [45] GenomicFeatures_1.38.0
## [46] backports_1.1.5
## [47] httpuv_1.5.2
## [48] Hmisc_4.2-0
```

```
## [49] tools_3.6.1
## [50] nor1mix_1.3-0
## [51] ggplot2_3.2.1
## [52] RColorBrewer_1.1-2
## [53] siggenes_1.60.0
## [54] Rcpp_1.0.2
## [55] plyr_1.8.4
## [56] base64enc_0.1-3
## [57] progress_1.2.2
## [58] zlibbioc_1.32.0
## [59] purrr_0.3.3
## [60] RCurl_1.95-4.12
## [61] BiasedUrn_1.07
## [62] prettyunits_1.0.2
## [63] rpart_4.1-15
## [64] openssl_1.4.1
## [65] cluster_2.1.0
## [66] magrittr_1.5
## [67] data.table_1.12.6
## [68] ProtGenerics_1.18.0
## [69] missMethyl_1.20.0
## [70] mime_0.7
## [71] hms_0.5.1
## [72] evaluate_0.14
## [73] xtable_1.8-4
## [74] XML_3.98-1.20
## [75] readxl_1.3.1
## [76] mclust_5.4.5
## [77] gridExtra_2.3
## [78] compiler_3.6.1
## [79] biomaRt_2.42.0
## [80] tibble_2.1.3
## [81] crayon_1.3.4
## [82] R.oo_1.22.0
## [83] htmltools_0.4.0
## [84] later_1.0.0
## [85] Formula_1.2-3
## [86] tidyr_1.0.0
## [87] DBI_1.0.0
## [88] MASS_7.3-51.4
## [89] rappdirs_0.3.1
## [90] Matrix_1.2-17
## [91] readr_1.3.1
## [92] permute_0.9-5
## [93] R.methodsS3_1.7.1
```

```
## [94] Gviz_1.30.0
## [95] pkgconfig_2.0.3
## [96] GenomicAlignments_1.22.0
## [97] registry_0.5-1
## [98] IlluminaHumanMethylation450kmanifest_0.4.0
## [99] foreign_0.8-72
## [100] xml2_1.2.2
## [101] annotate_1.64.0
## [102] rngtools_1.4
## [103] pkgmaker_0.27
## [104] multtest_2.42.0
## [105] beanplot_1.2
## [106] ruv_0.9.7.1
## [107] bibtex_0.4.2
## [108] doRNG_1.7.1
## [109] scrime_1.3.5
## [110] stringr_1.4.0
## [111] VariantAnnotation_1.32.0
## [112] digest_0.6.22
## [113] cellranger_1.1.0
## [114] base64_2.0
## [115] htmlTable_1.13.2
## [116] edgeR_3.28.0
## [117] DelayedMatrixStats_1.8.0
## [118] curl_4.2
## [119] shiny_1.4.0
## [120] Rsamtools_2.2.0
## [121] gtools_3.8.1
## [122] lifecycle_0.1.0
## [123] Rhdf5lib_1.8.0
## [124] askpass_1.1
## [125] limma_3.42.0
## [126] BSgenome_1.54.0
## [127] pillar_1.4.2
## [128] lattice_0.20-38
## [129] fastmap_1.0.1
## [130] httr_1.4.1
## [131] survival_2.44-1.1
## [132] GO.db_3.10.0
## [133] interactiveDisplayBase_1.24.0
## [134] glue_1.3.1
## [135] BiocVersion_3.10.1
## [136] bit_1.1-14
## [137] stringi_1.4.3
## [138] HDF5Array_1.14.0
```

```
## [139] blob_1.2.0
## [140] org.Hs.eg.db_3.10.0
## [141] latticeExtra_0.6-28
## [142] memoise_1.1.0
## [143] dplyr_0.8.3
```

References

- [1] Pidsley R, Zotenko E, Peters TJ, Lawrence MG, Risbridger GP, Molloy P, Van Dijk S, Muhlhäusler B, Stirzaker C, Clark SJ. Critical evaluation of the Illumina MethylationEPIC BeadChip microarray for whole-genome DNA methylation profiling. *Genome Biology*. 2016 17(1), 208.
- [2] Chen YA, Lemire M, Choufani S, Butcher DT, Grafodatskaya D, Zanke BW, Gallinger S, Hudson TJ, Weksberg R. Discovery of cross-reactive probes and polymorphic CpGs in the Illumina Infinium HumanMethylation450 microarray. *Epigenetics*. 2013 Jan 11;8(2).
- [3] Heyn H, Li N, Ferreira HJ, Moran S, Pisano DG, Gomez A, Esteller M. Distinct DNA methylomes of newborns and centenarians. *Proceedings of the National Academy of Sciences*. 2012 **109**(26), 10522-7.
- [4] Satterthwaite FE. An Approximate Distribution of Estimates of Variance Components., *Biometrics Bulletin*. 1946 **2**: 110-114
- [5] Chen Y, Pal B, Visvader JE, Smyth GK. Differential methylation analysis of reduced representation bisulfite sequencing experiments using edgeR. *F1000Research*. 2017 **6**, 2055.
- [6] Delacher M, Imbusch CD, Weichenhan D, Breiling A, Hotz-Wagenblatt A, Trager U, ... Feuerer M. (2017). Genome-wide DNA-methylation landscape defines specialization of regulatory T cells in tissues. *Nature Immunology*. 2017 **18**(10), 1160-1172.
- [7] Feng H, Conneely KN, Wu H. A Bayesian hierarchical model to detect differentially methylated loci from single nucleotide resolution sequencing data. *Nucleic Acids Research*. 2014 **42**(8), e69.
- [8] Park Y, and Wu H. Differential methylation analysis for BS-seq data under general experimental design. *Bioinformatics*. 2016 **32**(10), 1446-1453.