

An Introduction to the *REMP* Package

Yinan Zheng

September 25, 2017

Contents

1	Introduction	1
2	Installation	1
3	REMP: Repetitive Element Methylation Prediction	1
3.1	Groom methylation data	2
3.2	Prepare annotation data	2
3.3	Run prediction	3
3.4	Plot prediction	7
4	Session Information	8

1 Introduction

REMP predicts DNA methylation of locus-specific repetitive elements (RE) by learning surrounding genetic and epigenetic information. *REMP* provides genomewide single-base resolution of DNA methylation on RE that are difficult to measure using array-based or sequencing-based platforms, which enables epigenome-wide association study (EWAS) and differentially methylated region (DMR) analysis on RE.

2 Installation

Install *REMP* (release version):

```
> source("http://bioconductor.org/biocLite.R")
> biocLite("REMP")
```

To install devel version:

```
> library(devtools)
> install_github("YinanZheng/REMP")
```

Load *REMP* into the workspace:

```
> library(REMP)
```

3 REMP: Repetitive Element Methylation Prediction

Currently *REMP* supports Human (hg19, build 37) Alu and LINE-1 (L1) repetitive element (RE) methylation prediction using Illumina 450k or EPIC array.

3.1 Groom methylation data

Appropriate data preprocessing including quality control and normalization of methylation data are recommended before running *REMP*. Many packages are available to carry out these data preprocessing steps, for example, *minfi*, *wateRmelon*, and *methylumi*.

REMP is trying to minimize the requirement of the methylation data format. User can maintain the methylation data in *RatioSet* or *GenomicRatioSet* object offered by *minfi*, *data.table*, *data.frame*, *DataFrame*, or *matrix*. User can input either beta value or M-value. There are only two basic requirements of the methylation data:

1. Each row should represent CpG probe and each column should represent sample.
2. The row names should indicate Illumina probe ID (i.e. cg000000029).

However, there are some other common data issues that may prevent *REMP* from running correctly. For example, if the methylation data are in beta value and contain zero methylation values, logit transformation (to create M-value) will create negative infinite value; or the methylation data contain NA, Inf, or NaN data. To tackle these potential issues, *REMP* includes a handy function `groomMethy` which can help detect and fix these issues. We highly recommend to take advantage of this function:

```
> # Get GM12878 methylation data (450k array)
> GM12878_450k <- getGM12878('450k')
> GM12878_450k <- groomMethy(GM12878_450k, verbose = TRUE)
> GM12878_450k
```

```
class: RatioSet
dim: 482421 1
metadata(0):
assays(2): Beta M
rownames(482421): cg000000029 cg000000108 ...
               cg27666046 cg27666123
rowData names(0):
colnames(1): GM12878
colData names(0):
Annotation
  array: IlluminaHumanMethylation450k
  annotation: ilmn12.hg19
Preprocessing
  Method: NA
  minfi version: NA
  Manifest version: NA
```

For zero beta values, `groomMethy` will replace them with smallest non-zero beta value. For NA/NaN/Inf values, `groomMethy` will treat them as missing values and then apply KNN-imputation to complete the dataset. If the imputed value is out of the original range (which is possible when `imputebyrow = FALSE`), mean value will be used instead. Warning: imputed values for multimodal distributed CpGs (across samples) may not be correct. Please check package *ENmix* to identify the CpGs with multimodal distribution.

3.2 Prepare annotation data

To run *REMP* for RE methylation prediction, user first needs to prepare some annotation datasets. The function `initREMP` is designed to do the job.

Suppose user will predict Alu methylation using Illumina 450k array data:

```
> data(Alu.demo)
> remparcel <- initREMP(arrayType = "450k", REtype = "Alu",
+                       RE = Alu.demo, ncore = 1)
> remparcel
```

```
REMPParcel object
RE type: Alu
Illumina platform: 450k
Valid (max) RE-CpG flanking window size: 1200
Number of RE: 500
Number of RE-CpG: 5039
```

For demonstration, we only use 500 selected Alu sequence dataset which comes along with the package (`Alu.demo`). We specify `RE = Alu.demo`, so that the annotation dataset will be generated for the 500 selected Alu sequences. In real-world prediction, specifying RE is not necessary, as the function will pull up the complete RE sequence dataset from package *AnnotationHub*.

All data are stored in the *REMPParcel* object. It is recommended to specify a working directory so that the data generated can be preserved for later use:

```
> saveParcel(remparcel)
```

Without specifying working directory using option `work.dir`, the annotation dataset will be created under the temporal directory `tempdir()` by default. User can also turn on the `export` parameter in `initREMP` to save the data automatically.

3.3 Run prediction

Once the annotation data are ready, user can pass the annotation data parcel to `remp` for prediction:

```
> remp.res <- remp(GM12878_450k, REtype = 'Alu',
+                  parcel = remparcel, ncore = 1, seed = 777)
```

If `parcel` is missing, `remp` will then try to search the *REMPParcel* data file in the directory indicated by `work.dir` (use `tempdir()` if not specified).

By default, `remp` uses Random Forest (`method = 'rf'`) model (package *randomForest*) for prediction. Random Forest model is recommended because it offers more accurate prediction results and it automatically enables Quantile Regression Forest (Nicolai Meinshausen, 2006) for prediction reliability evaluation. `remp` constructs 19 predictors to carry out the prediction. For Random Forest model, the tuning parameter `param = 6` (i.e. `mtry` in *randomForest*) indicates how many predictors will be randomly selected for building the individual trees. The performance of random forest model is often relatively insensitive to the choice of `mtry`. Therefore, auto-tune will be turned off using random forest and `mtry` will be set to one third of the total number of predictors. It is recommended to specify a seed for reproducible prediction results.

`remp` will return a *REMPset* object, which inherits Bioconductor's *RangedSummarizedExperiment* class:

```
> remp.res

class: REMProduct
dim: 4808 1
metadata(8): REannotation RECPG ... GeneStats Seed
assays(3): rempB rempM rempQC
rownames: NULL
rowData names(1): RE.Index
colnames(1): GM12878
colData names(1): mtry

> # Display more detailed information
> details(remp.res)

RE type: Alu
Methylation profiling platform: 450k
Flanking window size: 1000
```

Prediction model: Random Forest
 QC model: Quantile Regression Forest
 Predicted 4808 CpG sites in 500 Alu

Number of predicted CpGs by chromosome:

chr1	chr2	chr3	chr4	chr5	chr6	chr7	chr8
449	276	293	131	179	397	292	102

chr9	chr10	chr11	chr12	chr13	chr14	chr15	chr16
98	148	254	310	66	127	133	333

chr17	chr18	chr19	chr20	chr21	chr22
295	81	674	66	37	67

Coverage information:

There are 500 profiled Alu by Illumina array.
 There are 481 RE-CpGs that have neighboring profiled CpGs are used for model training.
 In total, REMP predicts 500 Alu (4808 RE-CpG).
 Gene coverage by predicted Alu (out of total refSeq Gene):
 492 (1.96%) total genes;
 413 (2.15%) protein-coding genes;
 117 (1.59%) non-coding RNA genes.

Distribution of predicted methylation value (beta value):

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
0.02885	0.47313	0.66139	0.59494	0.75090	0.91733

Distribution of prediction reliability score:

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
0.6954	1.4627	1.7688	1.7957	2.0564	5.5713

Prediction results can be obtained by accessors:

```
> # Predicted RE-CpG methylation value (Beta value)
> rempB(remp.res)
```

DataFrame with 4808 rows and 1 column

```
GM12878
<numeric>
1 0.7649417
2 0.7883550
3 0.7967328
4 0.7991491
5 0.8018355
...
4804 0.4596627
4805 0.4585223
4806 0.4674303
4807 0.4888142
4808 0.4952238
```

```
> # Predicted RE-CpG methylation value (M value)
> rempM(remp.res)
```

DataFrame with 4808 rows and 1 column

```
GM12878
```

```

      <numeric>
1      1.702331
2      1.897199
3      1.970719
4      1.992340
5      2.016608
...      ...
4804 -0.23328479
4805 -0.23991011
4806 -0.18821933
4807 -0.06456144
4808 -0.02756350

```

```

> # Genomic location information of the predicted RE-CpG
> # Function inherit from class 'RangedSummarizedExperiment'
> rowRanges(remp.res)

```

GRanges object with 4808 ranges and 1 metadata column:

	seqnames	ranges	strand	RE.Index
	<Rle>	<IRanges>	<Rle>	<Rle>
[1]	chr1	[943927, 943928]	-	Alu_0000527
[2]	chr1	[943935, 943936]	-	Alu_0000527
[3]	chr1	[943968, 943969]	-	Alu_0000527
[4]	chr1	[943974, 943975]	-	Alu_0000527
[5]	chr1	[943991, 943992]	-	Alu_0000527
...
[4804]	chr22	[42095154, 42095155]	-	Alu_1170175
[4805]	chr22	[42095161, 42095162]	-	Alu_1170175
[4806]	chr22	[42095170, 42095171]	-	Alu_1170175
[4807]	chr22	[42095198, 42095199]	-	Alu_1170175
[4808]	chr22	[42095214, 42095215]	-	Alu_1170175

seqinfo: 93 sequences from an unspecified genome; no seqlengths

```

> # Standard error-scaled permutation importance of predictors
> imp(remp.res)

```

DataFrame with 19 rows and 1 column

```

      GM12878
      <numeric>
RE.score      7.844209
RE.Length     7.447086
RE.CpG.density 5.693255
RE.InNM       3.521871
RE.InNR       0.993378
...      ...
distance.min2 12.996170
Methy.min     28.932109
Methy.min2    11.232384
Methy.mean    12.948658
Methy.std     4.769137

```

```

> # Retrive seed number used for the reesults
> metadata(remp.res)$Seed

```

```

[1] 777

```

Trim off less reliable predicted results:

```
> # Any predicted CpG values with quality score < threshold (default = 1.7) will be replaced with NA. CpG
> # For mechanism study, more stringent cutoff is recommended.
> remp.res <- trim(remp.res)
> details(remp.res)
```

RE type: Alu

Methylation profiling platform: 450k

Flanking window size: 1000

Prediction model: Random Forest - trimmed (1.7)

QC model: Quantile Regression Forest

Predicted 2153 CpG sites in 392 Alu

Number of predicted CpGs by chromosome:

chr1	chr2	chr3	chr4	chr5	chr6	chr7	chr8
227	120	138	71	68	196	134	47

chr9	chr10	chr11	chr12	chr13	chr14	chr15	chr16
40	83	83	117	24	97	29	184

chr17	chr18	chr19	chr20	chr21	chr22
147	27	255	37	5	24

Coverage information:

There are 392 profiled Alu by Illumina array.

There are 355 RE-CpGs that have neighboring profiled CpGs are used for model training.

In total, REMP predicts 392 Alu (2153 RE-CpG).

Gene coverage by predicted Alu (out of total refSeq Gene):

375 (1.5%) total genes;

310 (1.62%) protein-coding genes;

87 (1.18%) non-coding RNA genes.

Distribution of predicted methylation value (beta value):

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
0.06761	0.67270	0.74426	0.70396	0.79209	0.91593

Distribution of prediction reliability score:

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
0.6954	1.2906	1.4331	1.4177	1.5700	1.7000

To add genomic regions annotation of the predicted REs:

```
> # By default gene symbol annotation will be added
> remp.res <- decodeAnnot(remp.res, ncore = 1)
> annotation(remp.res)
```

Seven genomic region indicators will be added to the annotation data in the input *REMPProduct* object:

- InNM: in protein-coding genes (overlap with refSeq gene's "NM" transcripts + 2000 bp upstream of the transcription start site (TSS))
- InNR: in noncoding RNA genes (overlap with refSeq gene's "NR" transcripts + 2000 bp upstream of the TSS)
- InTSS: in flanking region of 2000 bp upstream of the TSS. Default upstream limit is 2000 bp, which can be modified globally using `remp_options`

- In5UTR: in 5'untranslated regions (UTRs)
- InCDS: in coding DNA sequence regions
- InExon: in exon regions
- In3UTR: in 3'UTRs

Note that intron region and intergenic region information can be derived from the above genomic region indicators: if "InNM" and/or "InNR" is not missing but "InTSS", "In5UTR", "InExon", and "In3UTR" are missing, then the RE is strictly located within intron region; if all indicators are missing, then the RE is strictly located in intergenic region.

3.4 Plot prediction

Make a density plot of the predicted methylation (beta values):

```
> plot(remp.res, main = "Alu methylation (GM12878)", col = "blue")
```

4 Session Information

```
R version 3.4.1 (2017-06-30)
Platform: x86_64-w64-mingw32/x64 (64-bit)
Running under: Windows Server 2012 R2 x64 (build 9600)
```

```
Matrix products: default
```

```
locale:
```

```
[1] LC_COLLATE=C
[2] LC_CTYPE=English_United States.1252
[3] LC_MONETARY=English_United States.1252
[4] LC_NUMERIC=C
[5] LC_TIME=English_United States.1252
```

```
attached base packages:
```

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

```
other attached packages:
```

```
[1] REMP_1.0.7
[2] IlluminaHumanMethylationEPICanno.ilm10b2.hg19_0.6.0
[3] IlluminaHumanMethylation450kanno.ilmn12.hg19_0.6.0
[4] minfi_1.22.1
[5] bumphunter_1.16.0
[6] locfit_1.5-9.1
[7] iterators_1.0.8
[8] foreach_1.4.3
[9] Biostrings_2.44.2
[10] XVector_0.16.0
[11] SummarizedExperiment_1.6.5
[12] DelayedArray_0.2.7
[13] matrixStats_0.52.2
[14] Biobase_2.36.2
[15] GenomicRanges_1.28.5
[16] GenomeInfoDb_1.12.2
[17] IRanges_2.10.3
[18] S4Vectors_0.14.4
[19] BiocGenerics_0.22.0
[20] knitr_1.17
```

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

```
[1] colorspace_1.3-2
[2] class_7.3-14
[3] siggenes_1.50.0
[4] mclust_5.3
[5] base64_2.0
[6] DRR_0.0.2
[7] bit64_0.9-7
[8] interactiveDisplayBase_1.14.0
[9] AnnotationDbi_1.38.2
[10] prodlim_1.6.1
[11] lubridate_1.6.0
[12] ranger_0.8.0
```


[13] codetools_0.2-15
[14] splines_3.4.1
[15] doParallel_1.0.10
[16] impute_1.50.1
[17] robustbase_0.92-7
[18] RcppRoll_0.2.2
[19] Rsamtools_1.28.0
[20] caret_6.0-77
[21] annotate_1.54.0
[22] ddalpha_1.2.1
[23] kernlab_0.9-25
[24] settings_0.2.4
[25] shiny_1.0.5
[26] compiler_3.4.1
[27] httr_1.3.1
[28] assertthat_0.2.0
[29] Matrix_1.2-11
[30] lazyeval_0.2.0
[31] limma_3.32.7
[32] htmltools_0.3.6
[33] tools_3.4.1
[34] bindrcpp_0.2
[35] gtable_0.2.0
[36] glue_1.1.1
[37] GenomeInfoDbData_0.99.0
[38] reshape2_1.4.2
[39] dplyr_0.7.3
[40] doRNG_1.6.6
[41] Rcpp_0.12.12
[42] multtest_2.32.0
[43] preprocessCore_1.38.1
[44] nlme_3.1-131
[45] rtracklayer_1.36.4
[46] timeDate_3012.100
[47] gower_0.1.2
[48] stringr_1.2.0
[49] mime_0.5
[50] rngtools_1.2.4
[51] XML_3.98-1.9
[52] beanplot_1.2
[53] org.Hs.eg.db_3.4.1
[54] AnnotationHub_2.8.2
[55] DEoptimR_1.0-8
[56] zlibbioc_1.22.0
[57] MASS_7.3-47
[58] scales_0.5.0
[59] ipred_0.9-6
[60] BSgenome_1.44.2
[61] BiocInstaller_1.26.1
[62] GEOquery_2.42.0
[63] RColorBrewer_1.1-2
[64] curl_2.8.1
[65] yaml_2.1.14
[66] memoise_1.1.0

[67] ggplot2_2.2.1
[68] pkgmaker_0.22
[69] biomaRt_2.32.1
[70] rpart_4.1-11
[71] reshape_0.8.7
[72] stringi_1.1.5
[73] RSQLite_2.0
[74] genefilter_1.58.1
[75] GenomicFeatures_1.28.5
[76] BiocParallel_1.10.1
[77] lava_1.5
[78] rlang_0.1.2
[79] pkgconfig_2.0.1
[80] bitops_1.0-6
[81] nor1mix_1.2-3
[82] evaluate_0.10.1
[83] lattice_0.20-35
[84] purrr_0.2.3
[85] bindr_0.1
[86] GenomicAlignments_1.12.2
[87] recipes_0.1.0
[88] CVST_0.2-1
[89] bit_1.1-12
[90] BSgenome.Hsapiens.UCSC.hg19_1.4.0
[91] plyr_1.8.4
[92] magrittr_1.5
[93] R6_2.2.2
[94] dimRed_0.1.0
[95] DBI_0.7
[96] withr_2.0.0
[97] survival_2.41-3
[98] RCurl_1.95-4.8
[99] nnet_7.3-12
[100] tibble_1.3.4
[101] grid_3.4.1
[102] data.table_1.10.4
[103] blob_1.1.0
[104] ModelMetrics_1.1.0
[105] digest_0.6.12
[106] xtable_1.8-2
[107] httpuv_1.3.5
[108] illuminaio_0.18.0
[109] openssl_0.9.7
[110] munsell_0.4.3
[111] registry_0.3
[112] quadprog_1.5-5