# Introduction to RBM package

Dongmei Li

April 25, 2023

Clinical and Translational Science Institute, University of Rochester School of Medicine and Dentistry, Rochester, NY 14642-0708

## Contents

## 1 Overview

This document provides an introduction to the `RBM` package. The `RBM` package executes the resampling-based empirical Bayes approach using either permutation or bootstrap tests based on moderated t-statistics through the following steps.

- Firstly, the RBM package computes the moderated t-statistics based on the observed data set for each feature using the lmFit and eBayes function.

- Secondly, the original data are permuted or bootstrapped in a way that matches the null hypothesis to generate permuted or bootstrapped resamples, and the reference distribution is constructed using the resampled moderated t-statistics calculated from permutation or bootstrap resamples.

- Finally, the p-values from permutation or bootstrap tests are calculated based on the proportion of the permuted or bootstrapped moderated t-statistics that are as extreme as, or more extreme than, the observed moderated t-statistics.

Additional detailed information regarding resampling-based empirical Bayes approach can be found elsewhere (Li et al., 2013).

# 2 Getting started

The RBM package can be installed and loaded through the following R code.
Install the RBM package with:

```
> if (!requireNamespace("BiocManager", quietly=TRUE))
+     install.packages("BiocManager")
> BiocManager::install("RBM")
```

Load the RBM package with:

```
> library(RBM)
```

# 3 RBM_T and RBM_F functions

There are two functions in the RBM package: RBM_T and RBM_F. Both functions require input data in the matrix format with rows denoting features and columns denoting samples. RBM_T is used for two-group comparisons such as study designs with a treatment group and a control group. RBM_F can be used for more complex study designs such as more than two groups or time-course studies. Both functions need a vector for group notation, i.e., "1" denotes the treatment group and "0" denotes the control group. For the RBM_F function, a contrast vector need to be provided by users to perform pairwise comparisons between groups. For example, if the design has three groups (0, 1, 2), the aContrast parameter will be a vector such as ("X1-X0", "X2-X1", "X2-X0") to denote all pairwise comparisons. Users just need to add an extra "X" before the group labels to do the contrasts.

- Examples using the RBM_T function: normdata simulates a standardized gene expression data and unifdata simulates a methylation microarray data. The $p$-values from the RBM_T function could be further adjusted using the p.adjust function in the stats package through the Bejamini-Hochberg method.

```
> library(RBM)
> normdata <- matrix(rnorm(1000*6, 0, 1),1000,6)
> mydesign <- c(0,0,0,1,1,1)
> myresult <- RBM_T(normdata,mydesign,100,0.05)
> summary(myresult)

               Length Class  Mode
ordfit_t        1000   -none- numeric
ordfit_pvalue  1000   -none- numeric
ordfit_beta0   1000   -none- numeric
ordfit_beta1   1000   -none- numeric
permutation_p  1000   -none- numeric
bootstrap_p    1000   -none- numeric

> sum(myresult$permutation_p<=0.05)
```

2

```
[1] 21

> which(myresult$permutation_p<=0.05)

  [1]    9   40 196 211 272 328 343 349 352 391 447 493 641 649 668 704 748 773 816
[20] 859 945

> sum(myresult$bootstrap_p<=0.05)

[1] 7

> which(myresult$bootstrap_p<=0.05)

[1] 256 391 445 459 470 499 695

> permutation_adjp <- p.adjust(myresult$permutation_p, "BH")
> sum(permutation_adjp<=0.05)

[1] 2

> bootstrap_adjp <- p.adjust(myresult$bootstrap_p, "BH")
> sum(bootstrap_adjp<=0.05)

[1] 0

> unifdata <- matrix(runif(1000*7,0.10, 0.95), 1000, 7)
> mydesign2 <- c(0,0,0, 1,1,1,1)
> myresult2 <- RBM_T(unifdata,mydesign2,100,0.05)
> sum(myresult2$permutatioin_p<=0.05)

[1] 0

> sum(myresult2$bootstrap_p<=0.05)

[1] 10

> which(myresult2$bootstrap_p<=0.05)

  [1]   53   66 213 470 490 661 794 821 994 995

> bootstrap2_adjp <- p.adjust(myresult2$bootstrap_p, "BH")
> sum(bootstrap2_adjp<=0.05)

[1] 0
```

- Examples using the RBM_F function: normdata_F simulates a standardized gene expression data and unifdata_F simulates a methylation microarray data. In both examples, we were interested in pairwise comparisons.

3

```
> normdata_F <- matrix(rnorm(1000*9,0,2), 1000, 9)
> mydesign_F <- c(0, 0, 0, 1, 1, 1, 2, 2, 2)
> aContrast <- c("X1-X0", "X2-X1", "X2-X0")
> myresult_F <- RBM_F(normdata_F, mydesign_F, aContrast, 100, 0.05)
> summary(myresult_F)

               Length Class  Mode
ordfit_t       3000   -none- numeric
ordfit_pvalue  3000   -none- numeric
ordfit_beta1   3000   -none- numeric
permutation_p  3000   -none- numeric
bootstrap_p    3000   -none- numeric

> sum(myresult_F$permutation_p[, 1]<=0.05)

[1] 50

> sum(myresult_F$permutation_p[, 2]<=0.05)

[1] 43

> sum(myresult_F$permutation_p[, 3]<=0.05)

[1] 53

> which(myresult_F$permutation_p[, 1]<=0.05)

 [1]   37   66   94  122  123  163  166  167  170  223  235  263  277  292  294  325  340  364  368
[20]  411  414  463  467  485  489  492  515  549  599  621  631  671  674  677  682  685  714  725
[39]  734  780  789  794  860  870  879  912  942  958  977  994

> which(myresult_F$permutation_p[, 2]<=0.05)

 [1]   37   66   94  113  123  142  166  235  270  277  292  334  362  368  407  411  437  463  467
[20]  489  515  549  575  589  599  621  625  631  671  674  677  685  725  734  780  794  807  860
[39]  870  891  912  942  958

> which(myresult_F$permutation_p[, 3]<=0.05)

 [1]   37   39   66   94  113  116  122  142  187  223  277  292  324  364  367  368  376  403  411
[20]  463  467  489  507  515  549  570  575  589  599  621  625  631  654  671  674  677  682  685
[39]  734  774  780  794  807  809  837  860  869  870  879  891  912  958  994

> con1_adjp <- p.adjust(myresult_F$permutation_p[, 1], "BH")
> sum(con1_adjp<=0.05/3)

[1] 10
```

```
> con2_adjp <- p.adjust(myresult_F$permutation_p[, 2], "BH")
> sum(con2_adjp<=0.05/3)

[1] 5

> con3_adjp <- p.adjust(myresult_F$permutation_p[, 3], "BH")
> sum(con3_adjp<=0.05/3)

[1] 5

> which(con2_adjp<=0.05/3)

[1]   94 463 489 674 677

> which(con3_adjp<=0.05/3)

[1]   94 489 621 682 685

> unifdata_F <- matrix(runif(1000*18, 0.15, 0.98), 1000, 18)
> mydesign2_F <- c(rep(0, 6), rep(1, 6), rep(2, 6))
> aContrast <- c("X1-X0", "X2-X1", "X2-X0")
> myresult2_F <- RBM_F(unifdata_F, mydesign2_F, aContrast, 100, 0.05)
> summary(myresult2_F)

              Length Class  Mode
ordfit_t        3000   -none- numeric
ordfit_pvalue 3000   -none- numeric
ordfit_beta1  3000   -none- numeric
permutation_p 3000   -none- numeric
bootstrap_p   3000   -none- numeric

> sum(myresult2_F$bootstrap_p[, 1]<=0.05)

[1] 60

> sum(myresult2_F$bootstrap_p[, 2]<=0.05)

[1] 53

> sum(myresult2_F$bootstrap_p[, 3]<=0.05)

[1] 49

> which(myresult2_F$bootstrap_p[, 1]<=0.05)

 [1]   22   32   47   48   55   81   84   96  107  126  148  174  194  202  211  220  237  242  243
[20]  249  263  264  290  302  327  333  335  339  359  396  398  428  451  460  545  547  549  599
[39]  616  626  640  690  703  706  708  712  733  740  825  842  849  852  861  905  911  931  939
[58]  942  975  997
```

```
> which(myresult2_F$bootstrap_p[, 2]<=0.05)

 [1]   47   48   55   57   96  107  126  148  167  174  194  202  220  237  242  249  285  302  320
[20]  327  335  359  395  396  405  428  430  451  522  545  547  599  616  626  640  643  690  708
[39]  712  733  740  754  771  789  825  849  852  905  911  931  942  967  997

> which(myresult2_F$bootstrap_p[, 3]<=0.05)

 [1]   32   47   57   96  107  124  126  148  194  220  237  242  243  249  285  290  305  320  327
[20]  333  335  396  430  451  545  547  599  626  638  640  643  690  708  712  740  754  771  789
[39]  825  849  852  861  905  911  931  939  942  967  997

> con21_adjp <- p.adjust(myresult2_F$bootstrap_p[, 1], "BH")
> sum(con21_adjp<=0.05/3)

[1] 3

> con22_adjp <- p.adjust(myresult2_F$bootstrap_p[, 2], "BH")
> sum(con22_adjp<=0.05/3)

[1] 5

> con23_adjp <- p.adjust(myresult2_F$bootstrap_p[, 3], "BH")
> sum(con23_adjp<=0.05/3)

[1] 3
```

# 4 Ovarian cancer methylation example using the RBM_T function

Two-group comparisons are the most common contrast in biological and biomedical field. The ovarian cancer methylation example is used to illustrate the application of RBM_T in identifying differentially methylated loci. The ovarian cancer methylation example is taken from the gemone-wide DNA methylation profiling of United Kingdom Ovarian Cancer Population Study (UKOPS). This study used Illumina Infinium 27k Human DNA methylation Beadchip v1.2 to obtain DNA methylation profiles on over 27,000 CpGs in whole blood cells from 266 ovarian cancer women and 274 age-matched healthy controls. The data are downloaded from the NCBI GEO website with access number GSE19711. For illutration purpose, we chose the first 1000 loci in 8 randomly selected women with 4 ovariance cancer cases (pre-treatment) and 4 healthy controls. The following codes show the process of generating significant differential DNA methylation loci using the RBM_T function and presenting the results for further validation and investigations.

```
> system.file("data", package = "RBM")

[1] "/tmp/RtmpEdjKdh/Rinst17551bd43c6db/RBM/data"

> data(ovarian_cancer_methylation)
> summary(ovarian_cancer_methylation)
```

```
          IlmnID            Beta          exmdata2[, 2]       exmdata3[, 2]
 cg00000292:   1   Min.   :0.01058   Min.   :0.01187   Min.    :0.009103
 cg00002426:   1   1st Qu.:0.04111   1st Qu.:0.04407   1st Qu.:0.041543
 cg00003994:   1   Median :0.08284   Median :0.09531   Median :0.087042
 cg00005847:   1   Mean   :0.27397   Mean   :0.28872   Mean    :0.283729
 cg00006414:   1   3rd Qu.:0.52135   3rd Qu.:0.59032   3rd Qu.:0.558575
 cg00007981:   1   Max.   :0.97069   Max.   :0.96937   Max.    :0.970155
 (Other)   :994                      NA's   :4
 exmdata4[, 2]      exmdata5[, 2]     exmdata6[, 2]      exmdata7[, 2]
 Min.   :0.01019   Min.   :0.01108   Min.    :0.01937   Min.    :0.01278
 1st Qu.:0.04092   1st Qu.:0.04059   1st Qu.:0.05060   1st Qu.:0.04260
 Median :0.09042   Median :0.08527   Median :0.09502   Median :0.09362
 Mean   :0.28508   Mean   :0.28482   Mean    :0.27348   Mean    :0.27563
 3rd Qu.:0.57502   3rd Qu.:0.57300   3rd Qu.:0.52099   3rd Qu.:0.52240
 Max.   :0.96658   Max.   :0.97516   Max.    :0.96681   Max.    :0.95974
                   NA's   :1
 exmdata8[, 2]
 Min.   :0.01357
 1st Qu.:0.04387
 Median :0.09282
 Mean   :0.28679
 3rd Qu.:0.57217
 Max.   :0.96268

> ovarian_cancer_data <- ovarian_cancer_methylation[, -1]
> label <- c(1, 1, 0, 0, 1, 1, 0, 0)
> diff_results <- RBM_T(aData=ovarian_cancer_data, vec_trt=label, repetition=100, alpha=0.05)
> summary(diff_results)

               Length Class  Mode
ordfit_t        1000   -none- numeric
ordfit_pvalue   1000   -none- numeric
ordfit_beta0    1000   -none- numeric
ordfit_beta1    1000   -none- numeric
permutation_p   1000   -none- numeric
bootstrap_p     1000   -none- numeric

> sum(diff_results$ordfit_pvalue<=0.05)

[1] 45

> sum(diff_results$permutation_p<=0.05)

[1] 67

> sum(diff_results$bootstrap_p<=0.05)
```

```
[1] 53

> ordfit_adjp <- p.adjust(diff_results$ordfit_pvalue, "BH")
> sum(ordfit_adjp<=0.05)

[1] 0

> perm_adjp <- p.adjust(diff_results$permutation_p, "BH")
> sum(perm_adjp<=0.05)

[1] 5

> boot_adjp <- p.adjust(diff_results$bootstrap_p, "BH")
> sum(boot_adjp<=0.05)

[1] 0

> diff_list_perm <- which(perm_adjp<=0.05)
> diff_list_boot <- which(boot_adjp<=0.05)
> sig_results_perm <- cbind(ovarian_cancer_methylation[diff_list_perm, ], diff_results$ordfit_t[
> print(sig_results_perm)

        IlmnID       Beta exmdata2[, 2] exmdata3[, 2] exmdata4[, 2]
19  cg00016968 0.80628480            NA    0.81440820    0.83623180
245 cg00224508 0.04479948    0.04972043    0.04152814    0.04189373
450 cg00432979 0.03681359    0.04515700    0.04374394    0.03683598
627 cg00612467 0.04777553    0.03783457    0.05380982    0.05582291
764 cg00730260 0.90471270    0.90542290    0.91002680    0.91258610
    exmdata5[, 2] exmdata6[, 2] exmdata7[, 2] exmdata8[, 2]
19     0.80831380    0.73306440    0.82968340    0.84917800
245    0.04208405    0.05284988    0.03775905    0.03955271
450    0.04419125    0.04409653    0.02839263    0.03410020
627    0.04740551    0.05332965    0.05775211    0.05579710
764    0.90575890    0.88760470    0.90756300    0.90946790
    diff_results$ordfit_t[diff_list_perm]
19                           -2.446404
245                           1.962457
450                           1.546114
627                          -2.239498
764                          -1.808081
    diff_results$permutation_p[diff_list_perm]
19                                   0
245                                  0
450                                  0
627                                  0
764                                  0
```

```
> sig_results_boot <- cbind(ovarian_cancer_methylation[diff_list_boot, ], diff_results$ordfit_t[
> print(sig_results_boot)

 [1]  IlmnID
 [2]  Beta
 [3]  exmdata2[, 2]
 [4]  exmdata3[, 2]
 [5]  exmdata4[, 2]
 [6]  exmdata5[, 2]
 [7]  exmdata6[, 2]
 [8]  exmdata7[, 2]
 [9]  exmdata8[, 2]
[10]  diff_results$ordfit_t[diff_list_boot]
[11]  diff_results$bootstrap_p[diff_list_boot]
<0 rows> (or 0-length row.names)
```