# Introduction to RBM package

Dongmei Li

May 19, 2021

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)
```

```
[1] 16

> which(myresult$permutation_p<=0.05)

 [1]  39  49 119 122 155 195 270 314 322 388 392 426 733 785 799 993

> sum(myresult$bootstrap_p<=0.05)

[1] 3

> which(myresult$bootstrap_p<=0.05)

[1] 270 289 483

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

[1] 0

> 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] 29

> which(myresult2$bootstrap_p<=0.05)

 [1]  16  60 103 170 183 192 301 319 330 345 364 410 415 418 445 469 531 556 604
[20] 605 656 669 735 835 903 929 933 936 975

> 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.

```
> 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] 54

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

[1] 77

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

[1] 54

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

 [1]   17   36   48   92  102  113  119  130  147  149  173  188  193  215  216  233  236  258  289
[20]  306  343  351  361  363  377  397  406  415  478  541  561  591  596  598  609  690  710  719
[39]  726  728  736  783  789  807  826  833  861  875  878  925  931  932  960  987

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

 [1]   17   22   35   36   40   48   92  102  113  119  128  130  139  141  143  147  149  165  173
[20]  184  188  193  210  216  236  241  245  258  273  281  289  306  308  318  335  343  351  361
[39]  363  373  377  405  406  415  450  478  500  541  545  561  591  598  609  636  647  651  662
[58]  690  719  726  736  778  783  789  807  826  831  833  861  875  878  922  925  931  932  960
[77]  987

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

 [1]   17   36   48   56  102  113  119  128  130  141  147  149  173  188  193  215  216  236  245
[20]  258  289  306  335  343  351  361  363  377  406  415  478  519  541  588  591  598  609  655
[39]  690  726  736  748  789  807  826  833  861  875  878  909  925  931  932  987

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

```
[1] 11

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

[1] 18

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

[1] 11

> which(con2_adjp<=0.05/3)

 [1]   40   48  113  119  289  351  361  415  541  598  636  736  789  807  826  875  878  987

> which(con3_adjp<=0.05/3)

 [1]   36   48  113  343  541  690  736  807  826  925  987

> 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] 58

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

[1] 51

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

[1] 57

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

```
   [1]   43   48   49   56   75   83   93   97   98  107  115  127  132  134  135  137  138  147  180
  [20]  197  199  223  275  298  301  313  329  384  400  427  461  471  473  494  535  563  714  736
  [39]  756  780  787  806  813  876  880  889  911  912  926  928  929  946  952  963  974  986  993
  [58]  999

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

   [1]   43   48   75   76   83   93   97  107  115  127  132  135  137  147  165  197  213  275  298
  [20]  313  329  383  390  427  445  471  473  494  520  535  710  714  736  756  771  778  780  806
  [39]  813  876  879  911  912  917  926  928  929  946  952  974  986

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

   [1]   43   48   75   83   93   97  107  134  135  137  138  147  149  165  180  183  197  199  223
  [20]  275  298  301  313  329  364  383  384  390  401  427  445  473  494  520  535  710  714  736
  [39]  780  787  806  813  876  880  889  911  912  917  926  928  929  946  952  963  974  986  999

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

[1] 9

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

[1] 7

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

[1] 10
```

# 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] "/private/tmp/RtmpPcJupN/Rinstfc3274b1ebb3/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] 47

> sum(diff_results$bootstrap_p<=0.05)

[1] 42

> 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$ordfit_t
> print(sig_results_perm)
        IlmnID      Beta exmdata2[, 2] exmdata3[, 2] exmdata4[, 2]
5   cg00006414 0.07635468    0.07442468    0.15698040    0.08676092
245 cg00224508 0.04479948    0.04972043    0.04152814    0.04189373
346 cg00331237 0.05972383            NA    0.08204769    0.08345662
627 cg00612467 0.04777553    0.03783457    0.05380982    0.05582291
764 cg00730260 0.90471270    0.90542290    0.91002680    0.91258610
    exmdata5[, 2] exmdata6[, 2] exmdata7[, 2] exmdata8[, 2]
5      0.07982556    0.08111396    0.08271889    0.08045977
245    0.04208405    0.05284988    0.03775905    0.03955271
346    0.05372019    0.06241126    0.06955040    0.09140985
627    0.04740551    0.05332965    0.05775211    0.05579710
764    0.90575890    0.88760470    0.90756300    0.90946790
    diff_results$ordfit_t[diff_list_perm]
5                                -1.389459
245                               1.962457
346                              -3.767916
627                              -2.239498
764                              -1.808081
    diff_results$permutation_p[diff_list_perm]
5                                            0
245                                          0
346                                          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)
```