# Introduction to RBM package

Dongmei Li

October 21, 2024

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] 61

> which(myresult$permutation_p<=0.05)

 [1]   53   57   66   80  100  125  134  146  198  204  205  235  263  287  303  339  348  363  364
[20]  373  407  441  470  472  486  496  519  528  590  593  616  618  629  643  667  705  718  734
[39]  746  760  775  782  789  794  818  821  822  836  875  904  918  922  923  929  931  938  950
[58]  953  986  992  998

> sum(myresult$bootstrap_p<=0.05)

[1] 4

> which(myresult$bootstrap_p<=0.05)

[1] 734 746 756 904

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

[1] 4

> 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] 15

> which(myresult2$bootstrap_p<=0.05)

 [1]   34   43   90  108  148  201  267  332  334  383  578  616  799  887  988

> 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] 53

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

[1] 60

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

[1] 52

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

 [1]   42   57   74   95  117  161  182  187  204  234  239  250  258  279  292  294  301  306  311
[20]  318  328  382  408  440  447  479  485  498  506  531  540  583  595  630  633  653  655  694
[39]  702  745  759  783  815  819  820  835  854  879  898  921  950  968  995

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

 [1]    3   42   57   74   94   95  117  141  161  182  187  234  239  250  253  258  292  294  301
[20]  306  311  318  328  341  382  408  447  479  485  498  506  531  532  540  583  595  602  625
[39]  633  646  651  653  655  678  693  694  702  743  745  759  783  815  819  835  854  879  896
[58]  921  924  950

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

 [1]    3   42   57   70   95  117  141  161  182  187  234  239  250  269  279  294  301  306  311
[20]  318  328  335  382  408  479  485  506  531  532  540  583  602  625  633  653  678  702  745
[39]  759  783  794  815  820  854  879  896  898  921  950  966  991  995
```

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

[1] 18

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

[1] 11

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

[1] 2

> which(con2_adjp<=0.05/3)

 [1] 187 250 301 328 479 506 655 702 759 783 879

> which(con3_adjp<=0.05/3)

[1] 187 653

> 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] 52

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

[1] 51

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

[1] 59
```

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

 [1]    5    8   21   36  119  124  180  181  191  235  246  259  269  277  282  290  304  333  371
[20]  391  431  455  472  474  495  498  501  522  549  566  607  644  677  716  724  745  747  748
[39]  786  787  796  807  811  840  859  860  906  923  927  937  956  962

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

 [1]    8   42   52   60  119  124  180  181  191  211  235  269  275  277  282  290  325  333  337
[20]  371  455  472  474  495  498  501  515  517  522  537  644  724  738  745  747  748  787  796
[39]  807  811  819  834  840  843  860  886  906  909  923  956  995

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

 [1]    5    8   21   36   73  119  124  180  181  191  211  231  235  275  277  282  290  333  337
[20]  371  391  431  453  455  472  474  495  498  501  522  537  540  566  644  672  716  722  724
[39]  738  745  747  748  749  786  787  794  796  807  811  834  840  860  862  879  906  909  923
[58]  956  992

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

[1] 0

> 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] 5
```

# 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] "F:/biocbuild/bbs-3.20-bioc/tmpdir/RtmpEhlHKw/Rinst164742db64780/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] 47
```

7

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

[1] 53

> sum(diff_results$bootstrap_p<=0.05)

[1] 44

> 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] 11

> 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
97  cg00083937 0.53046980    0.60529020    0.62733150    0.65623920
103 cg00094319 0.73784280    0.73532960    0.75574900    0.73830220
131 cg00121904 0.15449580    0.17949750    0.23608110    0.24354150
259 cg00234961 0.04192170    0.04321576    0.05707140    0.05327565
285 cg00263760 0.09050395    0.10197760    0.14801710    0.12242400
627 cg00612467 0.04777553    0.03783457    0.05380982    0.05582291
764 cg00730260 0.90471270    0.90542290    0.91002680    0.91258610
851 cg00830029 0.58362500    0.59397870    0.64739610    0.67269640
887 cg00862290 0.43640520    0.54047160    0.60786800    0.56325950
928 cg00901493 0.03737166    0.03903724    0.04684618    0.04981432
    exmdata5[, 2] exmdata6[, 2] exmdata7[, 2] exmdata8[, 2]
19     0.80831380    0.73306440    0.82968340    0.84917800
97     0.55974270    0.43157020    0.64046990    0.57876990
103    0.67349260    0.73510200    0.75715920    0.78981220
131    0.17352980    0.12564280    0.18193170    0.20847670
259    0.04030003    0.03996053    0.05086962    0.05445672
285    0.11693600    0.10650430    0.12281160    0.12310430
```

```
627     0.04740551      0.05332965      0.05775211      0.05579710
764     0.90575890      0.88760470      0.90756300      0.90946790
851     0.50820240      0.34657470      0.66276570      0.64634510
887     0.50259740      0.40111730      0.56646700      0.54552980
928     0.04490690      0.04204062      0.05050039      0.05268215
    diff_results$ordfit_t[diff_list_perm]
19                              -2.547097
97                              -2.665377
103                             -2.343784
131                             -3.562745
259                             -2.833203
285                             -2.993292
627                             -1.797392
764                             -1.560713
851                             -2.986319
887                             -3.368752
928                             -1.982308
    diff_results$permutation_p[diff_list_perm]
19                                       0
97                                       0
103                                      0
131                                      0
259                                      0
285                                      0
627                                      0
764                                      0
851                                      0
887                                      0
928                                      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)
```