# Introduction to RBM package

Dongmei Li

November 1, 2022

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] 36

> which(myresult$permutation_p<=0.05)

 [1]   13   25  160  234  249  277  297  334  335  337  360  375  435  445  447  479  483  541  547
[20] 549  554  565  593  614  630  659  663  674  718  723  732  797  881  899  955  960

> sum(myresult$bootstrap_p<=0.05)

[1] 7

> which(myresult$bootstrap_p<=0.05)

[1] 232 391 445 455 549 960 986

> 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] 29

> which(myresult2$bootstrap_p<=0.05)

 [1]   31   34   93  138  184  202  210  249  361  374  377  386  401  462  470  471  485  527  563
[20] 569  575  599  656  670  704  785  790  860  982

> 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] 59

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

[1] 76

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

[1] 75

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

 [1]    5   37   39   45   50   66   81   83   95  105  109  127  134  157  159  173  184  187  190
[20]  203  244  248  254  261  266  270  282  293  295  345  349  360  381  409  416  488  510  533
[39]  570  576  587  596  608  678  696  699  739  777  879  894  901  902  923  940  962  971  973
[58]  980  981

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

 [1]    5   37   39   45   50   66   81   83   95  105  157  159  167  184  187  190  193  203  215
[20]  244  248  254  261  266  270  274  282  293  295  333  345  348  349  360  364  381  409  416
[39]  426  466  488  533  551  570  576  587  596  608  613  624  666  678  696  699  702  739  745
[58]  769  777  790  795  864  879  888  894  901  902  923  929  940  962  971  973  980  981  995

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

 [1]    5   37   39   45   50   66   81   83   95  105  115  117  127  134  144  157  159  163  167
[20]  173  184  187  190  195  198  215  244  248  261  266  270  282  293  295  340  345  348  349
[39]  360  364  381  409  416  466  488  533  551  570  576  579  587  596  613  619  678  690  696
[58]  699  739  777  790  795  879  888  894  896  901  902  923  929  940  962  971  973  980

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

4

```
[1] 6

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

[1] 12

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

[1] 14

> which(con2_adjp<=0.05/3)

 [1]   45   83 266 270 282 293 345 533 596 790 879 902

> which(con3_adjp<=0.05/3)

 [1]    5   39   45   66   95 105 270 293 416 587 596 902 962 980

> 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] 54

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

[1] 42

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

[1] 52

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

```
 [1]  13   42   68   89  110  119  142  167  194  207  238  264  265  276  290  291  313  338  367
[20] 387  389  402  443  446  450  454  458  472  479  512  523  527  528  529  533  538  598  613
[39] 619  624  635  636  653  726  727  731  790  833  890  925  935  942  947  961

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

 [1]  13   42   50   68  119  130  142  207  208  238  264  265  276  291  313  387  389  402  443
[20] 450  458  479  512  528  529  617  624  635  636  653  668  726  727  767  777  790  890  935
[39] 942  947  961  973

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

 [1]  13   50   55   59   68   89  111  119  121  130  142  167  194  207  208  238  264  265  276
[20] 288  291  313  343  387  389  443  450  458  462  479  504  512  528  533  538  624  635  636
[39] 653  726  727  731  767  788  790  890  933  935  942  947  961  973

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

[1] 7

> 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] 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 gemonewide 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/var/folders/db/4tvgx8jx4z3fm1gzlnlzw9rc0000gq/T/RtmpqrmQ5A/Rinst103b5792f86d7/RBM/

> 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] 62

> sum(diff_results$bootstrap_p<=0.05)

[1] 66

> 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] 7

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

[1] 10

> 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]
16   cg00014085 0.05906804    0.04518973    0.04211710    0.03665208
83   cg00072216 0.04505377    0.04598964    0.04000674    0.03231534
237  cg00215066 0.94926640    0.95311870    0.94634910    0.94561120
245  cg00224508 0.04479948    0.04972043    0.04152814    0.04189373
280  cg00260778 0.64319890    0.60488960    0.56735060    0.53150910
437  cg00424946 0.04122172    0.04325330    0.03339863    0.02876798
931  cg00901704 0.05734342    0.04812868    0.04478214    0.03878488
    exmdata5[, 2] exmdata6[, 2] exmdata7[, 2] exmdata8[, 2]
16     0.04222944    0.05324246    0.03728026    0.04062589
83     0.04965089    0.04833366    0.03466159    0.04390894
237    0.94837410    0.94665570    0.94089070    0.94600090
245    0.04208405    0.05284988    0.03775905    0.03955271
280    0.61920530    0.61925200    0.46753250    0.55632410
437    0.03353116    0.03719167    0.03096761    0.03234779
931    0.04497277    0.05751033    0.03089829    0.04423603
    diff_results$ordfit_t[diff_list_perm]
16                           2.325659
83                           2.514109
237                          1.419654
245                          1.962457
```

```
280                                          4.170347
437                                          2.102892
931                                          2.464709
     diff_results$permutation_p[diff_list_perm]
16                                                  0
83                                                  0
237                                                 0
245                                                 0
280                                                 0
437                                                 0
931                                                 0


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

         IlmnID        Beta exmdata2[, 2] exmdata3[, 2] exmdata4[, 2]
95   cg00081975 0.03633894    0.04975194    0.06024723    0.05598723
106  cg00095674 0.07076291    0.05045181    0.03861991    0.03337576
131  cg00121904 0.15449580    0.17949750    0.23608110    0.24354150
146  cg00134539 0.61101320    0.53321780    0.45999340    0.46787420
259  cg00234961 0.04192170    0.04321576    0.05707140    0.05327565
280  cg00260778 0.64319890    0.60488960    0.56735060    0.53150910
350  cg00333226 0.08320694    0.11587190    0.14999650    0.14307800
911  cg00888479 0.07388961    0.07361080    0.10149800    0.09985076
928  cg00901493 0.03737166    0.03903724    0.04684618    0.04981432
979  cg00945507 0.13432250    0.23854600    0.34749760    0.28903340
     exmdata5[, 2] exmdata6[, 2] exmdata7[, 2] exmdata8[, 2]
95      0.04561792    0.05115624    0.06068253    0.06168212
106     0.04693030    0.06837343    0.04534005    0.03709488
131     0.17352980    0.12564280    0.18193170    0.20847670
146     0.67191510    0.63137380    0.47929610    0.45428300
259     0.04030003    0.03996053    0.05086962    0.05445672
280     0.61920530    0.61925200    0.46753250    0.55632410
350     0.10704480    0.13751630    0.12588230    0.13863730
911     0.08633986    0.06765189    0.09070268    0.12417730
928     0.04490690    0.04204062    0.05050039    0.05268215
979     0.11848510    0.16653850    0.30718420    0.26624740
     diff_results$ordfit_t[diff_list_boot]
95                                -3.252063
106                                3.100324
131                               -3.451679
146                                5.394750
259                               -4.052697
280                                4.170347
350                               -2.458696
911                               -3.621731
```

```
928                              -2.716443
979                              -4.750997
    diff_results$bootstrap_p[diff_list_boot]
95                                      0
106                                     0
131                                     0
146                                     0
259                                     0
280                                     0
350                                     0
911                                     0
928                                     0
979                                     0
```