# Introduction to RBM package

Dongmei Li

May 1, 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)
```

```
[1] 47

> which(myresult$permutation_p<=0.05)

 [1]    6   28   81   87  108  131  172  187  191  209  232  263  276  287  305  335  348  350  362
[20]  370  373  377  380  399  423  447  462  485  567  586  594  599  600  613  639  666  703  758
[39]  834  843  927  954  955  960  977  982  998

> sum(myresult$bootstrap_p<=0.05)

[1] 7

> which(myresult$bootstrap_p<=0.05)

[1] 131 192 215 492 650 720 923

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

[1] 12

> 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] 27

> which(myresult2$bootstrap_p<=0.05)

 [1]   21   85   87   94  103  175  236  269  320  333  402  436  464  501  537  566  593  612  666
[20]  743  762  817  839  849  863  935  985

> 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] 52

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

[1] 41

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

[1] 49

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

 [1]   29   46   57   62   63   76   86  139  142  167  171  172  183  219  332  364  417  437  460
[20]  475  491  503  514  539  548  577  617  620  622  655  667  705  708  715  734  762  806  811
[39]  821  822  824  829  839  840  874  875  882  885  922  932  964  976

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

 [1]   29   51   62   63   76   86   90  139  142  172  183  219  313  332  363  402  470  482  491
[20]  503  514  539  577  617  620  633  655  659  667  705  762  806  810  822  824  829  840  875
[39]  882  948  964

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

 [1]   29   62   63   76   84   86   90  139  142  171  172  183  219  313  332  364  402  437  470
[20]  503  507  514  532  539  549  577  617  620  622  633  655  667  705  715  762  806  810  811
[39]  822  824  829  840  874  875  882  885  922  964  976
```

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

[1] 9

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

[1] 8

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

[1] 7

> which(con2_adjp<=0.05/3)

[1]   29   76   86 514 667 705 762 875

> which(con3_adjp<=0.05/3)

[1]   86 142 514 620 667 762 875

> 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] 61

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

[1] 59

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

[1] 54
```

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

 [1]   28   35   41   73  125  129  156  197  203  205  218  221  238  272  276  289  305  310  335
[20]  360  397  416  436  447  449  460  461  489  495  501  510  536  557  628  635  637  650  657
[39]  662  700  704  714  740  752  755  757  789  792  799  819  832  837  861  876  882  883  903
[58]  915  976  978  988

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

 [1]   28   41   43   55   73   80  125  129  156  177  197  203  205  208  218  221  233  238  272
[20]  276  289  305  335  397  416  436  447  449  460  461  474  489  501  510  535  536  539  575
[39]  615  628  635  637  657  658  662  700  702  755  757  819  837  861  876  882  883  896  903
[58]  915  978

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

 [1]   28   35   41  125  129  156  177  197  203  208  218  221  233  266  272  276  284  289  305
[20]  335  358  360  436  449  460  489  501  510  536  557  628  635  637  650  657  700  702  704
[39]  714  740  755  757  789  792  819  832  837  876  882  883  903  915  978  988

> 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] 4

> 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/RtmpLXoC6n/Rinstfc4dc34b50e1a/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] 57

> sum(diff_results$bootstrap_p<=0.05)

[1] 72

> 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] 3

> 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$ordfit_t[
> print(sig_results_perm)
        IlmnID       Beta exmdata2[, 2] exmdata3[, 2] exmdata4[, 2]
245 cg00224508 0.04479948    0.04972043    0.04152814    0.04189373
280 cg00260778 0.64319890    0.60488960    0.56735060    0.53150910
764 cg00730260 0.90471270    0.90542290    0.91002680    0.91258610
    exmdata5[, 2] exmdata6[, 2] exmdata7[, 2] exmdata8[, 2]
245    0.04208405    0.05284988    0.03775905    0.03955271
280    0.61920530    0.61925200    0.46753250    0.55632410
764    0.90575890    0.88760470    0.90756300    0.90946790
    diff_results$ordfit_t[diff_list_perm]
245                              1.962457
280                              4.170347
764                             -1.808081
    diff_results$permutation_p[diff_list_perm]
245                                          0
280                                          0
764                                          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]
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
346 cg00331237 0.05972383          NA      0.08204769    0.08345662
397 cg00394658 0.27940900    0.40410330    0.40262320    0.44339290
677 cg00651216 0.06825629    0.12529090    0.14409190    0.13907250
833 cg00814580 0.09348613    0.09619816    0.12010440    0.11534240
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]
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
346    0.05372019    0.06241126    0.06955040    0.09140985
397    0.35626060    0.23388380    0.41974630    0.45806880
677    0.07669587    0.09597587    0.11690440    0.15194540
833    0.09577040    0.11598850    0.12860890    0.14111200
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$ordfit_t[diff_list_boot]
146                          5.394750
259                         -4.052697
280                          4.170347
346                         -3.767916
397                         -3.070559
677                         -3.387628
833                         -3.428319
911                         -3.621731
928                         -2.716443
979                         -4.750997
    diff_results$bootstrap_p[diff_list_boot]
146                                    0
259                                    0
280                                    0
346                                    0
397                                    0
677                                    0
833                                    0
911                                    0
928                                    0
979                                    0
```