

# Introduction to RBM package

Dongmei Li

May 3, 2016

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:

```
> source("http://bioconductor.org/biocLite.R")
> biocLite("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] 9

> which(myresult$permutation_p<=0.05)
[1] 111 155 307 436 452 668 777 808 850

> sum(myresult$bootstrap_p<=0.05)
[1] 7

> which(myresult$bootstrap_p<=0.05)
[1] 162 308 497 575 668 808 966
```

```

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

[1] 1

> 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] 0

> which(myresult2$bootstrap_p<=0.05)

integer(0)

> 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] 64

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

[1] 70

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

[1] 78

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

[1]   6   8  12  30  58  64  82  83  86  98  99 156 183 192 221 242 279 321 326
[20] 330 332 335 355 370 374 386 396 399 440 448 481 489 521 545 554 591 599 611
[39] 621 642 648 655 664 674 747 761 771 773 779 784 815 816 837 842 864 886 894
[58] 900 901 907 931 956 984 996

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

[1]   6   8  12  21  30  58  64  82  83  86  98  156 183 192 214 217 221 230 242
[20] 268 279 326 330 332 335 355 358 370 374 386 396 399 440 448 461 481 488 489
[39] 513 521 554 562 568 591 599 611 621 642 648 655 664 673 674 747 771 773 779
[58] 815 816 837 842 864 886 894 900 901 903 983 984 996

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

[1]   3   6   8  12  30  58  64  82  86  98  99 156 183 192 217 218 221 238 242
[20] 255 268 326 330 332 335 355 358 370 374 380 386 396 399 411 440 448 461 481
[39] 489 513 521 532 545 554 562 569 591 599 611 621 642 645 648 655 664 674 690
[58] 747 771 773 779 784 815 816 837 842 864 880 886 894 896 900 901 903 931 983
[77] 984 996

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

[1] 13

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

[1] 10

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

[1] 19

> which(con2_adjp<=0.05/3)

```

```

[1] 30 64 374 399 448 521 771 815 816 886
> which(con3_adjp<=0.05/3)
[1] 8 12 30 64 86 192 326 374 481 521 655 664 747 771 815 842 886 894 901

> 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] 45

> sum(myresult2_F$bootstrap_p[, 2]<=0.05)
[1] 52

> sum(myresult2_F$bootstrap_p[, 3]<=0.05)
[1] 41

> which(myresult2_F$bootstrap_p[, 1]<=0.05)
[1] 3 64 83 137 170 183 188 202 284 301 311 358 359 375 391 420 434 435 447
[20] 467 534 540 543 554 565 593 604 618 629 662 700 708 753 764 772 776 817 844
[39] 852 866 881 959 979 986 995

> which(myresult2_F$bootstrap_p[, 2]<=0.05)
[1] 3 12 64 83 100 137 170 183 188 202 231 301 311 358 359 375 391 420 431
[20] 434 435 447 467 473 496 534 536 554 565 566 593 618 621 629 643 662 700 708
[39] 764 772 776 817 840 844 852 881 900 939 959 972 986 995

> which(myresult2_F$bootstrap_p[, 3]<=0.05)
[1] 3 64 83 137 170 183 188 202 284 301 311 358 359 375 391 420 431 434 447
[20] 448 467 496 534 554 565 618 629 700 708 772 776 790 817 840 844 852 866 881
[39] 972 986 995

```

```

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

[1] 6

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

[1] 8

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

[1] 7

```

## 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 genome-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 illustration purpose, we chose the first 1000 loci in 8 randomly selected women with 4 ovarian 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/Rtmp6t1opA/Rinst37ff1739512d/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] 59

> 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] 4

```

```

> 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[, diff_list_perm])
> print(sig_results_perm)

    IlmnID      Beta exmdata2[, 2] exmdata3[, 2] exmdata4[, 2]
245 cg00224508 0.04479948   0.04972043   0.04152814   0.04189373
627 cg00612467 0.04777553   0.03783457   0.05380982   0.05582291
764 cg00730260 0.90471270   0.90542290   0.91002680   0.91258610
928 cg00901493 0.03737166   0.03903724   0.04684618   0.04981432
               exmdata5[, 2] exmdata6[, 2] exmdata7[, 2] exmdata8[, 2]
245     0.04208405   0.05284988   0.03775905   0.03955271
627     0.04740551   0.05332965   0.05775211   0.05579710
764     0.90575890   0.88760470   0.90756300   0.90946790
928     0.04490690   0.04204062   0.05050039   0.05268215
    diff_results$ordfit_t[diff_list_perm]
245                      1.962457
627                     -2.239498
764                     -1.808081
928                     -2.716443
    diff_results$permutation_p[diff_list_perm]
245                         0
627                         0
764                         0
928                         0

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

    IlmnID      Beta exmdata2[, 2] exmdata3[, 2] exmdata4[, 2]
16 cg00014085 0.05906804   0.04518973   0.04211710   0.03665208
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
285 cg00263760 0.09050395   0.10197760   0.14801710   0.12242400
743 cg00717862 0.07999436   0.07873347   0.06089359   0.06171374
851 cg00830029 0.58362500   0.59397870   0.64739610   0.67269640
887 cg00862290 0.43640520   0.54047160   0.60786800   0.56325950
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

```

```

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
285 0.11693600 0.10650430 0.12281160 0.12310430
743 0.07594936 0.09062161 0.06475791 0.07271878
851 0.50820240 0.34657470 0.66276570 0.64634510
887 0.50259740 0.40111730 0.56646700 0.54552980
931 0.04497277 0.05751033 0.03089829 0.04423603

diff_results$ordfit_t[diff_list_boot]
16 2.325659
131 -3.451679
146 5.394750
259 -4.052697
280 4.170347
285 -3.093997
743 3.444684
851 -2.841244
887 -3.217939
931 2.464709

diff_results$bootstrap_p[diff_list_boot]
16 0
131 0
146 0
259 0
280 0
285 0
743 0
851 0
887 0
931 0

```