

# Introduction to RBM package

Dongmei Li

April 30, 2018

Clinical and Translational Science Institute, University of Rochester School of Medicine and Dentistry, Rochester, NY 14642-0708

## Contents

<b>1 Overview</b>	<b>1</b>
<b>2 Getting started</b>	<b>2</b>
<b>3 RBM_T and RBM_F functions</b>	<b>2</b>
<b>4 Ovarian cancer methylation example using the RBM_T function</b>	<b>6</b>

## 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] 52
```

```

> which(myresult$permutation_p<=0.05)

[1] 2 25 29 71 93 104 112 138 173 239 269 294 297 300 323 328 333 348 357
[20] 373 411 413 416 449 454 471 487 497 512 527 546 570 598 639 640 653 691 692
[39] 705 713 728 744 749 767 796 812 848 901 923 966 968 973

> sum(myresult$bootstrap_p<=0.05)

[1] 8

> which(myresult$bootstrap_p<=0.05)

[1] 64 354 357 411 527 705 715 975

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

[1] 8

> 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$permutation_p<=0.05)

[1] 0

> sum(myresult2$bootstrap_p<=0.05)

[1] 6

> which(myresult2$bootstrap_p<=0.05)

[1] 34 347 618 740 781 900

> 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] 56

> sum(myresult_F$permutation_p[, 2]<=0.05)
[1] 63

> sum(myresult_F$permutation_p[, 3]<=0.05)
[1] 55

> which(myresult_F$permutation_p[, 1]<=0.05)
[1]  24  25  75  76  95 106 124 136 139 161 174 258 280 294 296 303 350 351 353
[20] 392 398 416 423 439 451 452 464 474 505 514 518 536 537 560 569 577 580 614
[39] 622 642 644 653 661 706 748 758 788 794 801 812 846 856 907 975 984 988

> which(myresult_F$permutation_p[, 2]<=0.05)
[1]  24  25  69  75  84  95 124 136 139 140 173 214 294 296 303 350 353 392 398
[20] 403 416 422 439 451 452 464 505 514 518 523 536 537 560 569 580 603 606 614
[39] 622 642 653 661 685 695 706 714 716 748 758 788 791 794 801 812 827 835 846
[58] 856 861 889 907 984 988

> which(myresult_F$permutation_p[, 3]<=0.05)
[1]  24  75  95 110 124 139 173 174 234 258 294 296 303 350 353 392 403 416 421
[20] 423 439 451 452 464 505 514 536 537 560 564 569 577 579 580 603 614 622 642
[39] 653 661 685 706 716 748 762 788 791 794 801 812 846 856 975 984 988

> 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] 9

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

[1] 10

> which(con2_adjp<=0.05/3)

[1] 95 353 451 536 788 794 801 812 984

> which(con3_adjp<=0.05/3)

[1] 24 353 416 451 536 794 801 812 846 984

> 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] 57

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

[1] 55

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

[1] 56

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

[1] 26 89 125 146 167 170 174 191 211 219 238 285 290 306 326 351 404 409 415
[20] 432 435 441 451 546 559 561 580 598 612 619 626 635 649 671 701 720 786 792
[39] 795 824 832 834 856 866 879 888 891 907 909 939 949 956 962 969 971 994 999

```

```

> which(myresult2_F$bootstrap_p[, 2]<=0.05)
[1] 26 38 89 125 128 145 146 167 170 191 211 219 238 248 285 290 306 313 326
[20] 330 409 415 432 441 451 559 561 585 598 612 626 635 649 658 680 694 720 786
[39] 795 824 832 834 840 856 866 879 907 909 939 949 962 969 971 973 999

> which(myresult2_F$bootstrap_p[, 3]<=0.05)
[1] 26 58 61 125 146 167 170 174 191 211 219 238 285 290 306 326 404 409 415
[20] 432 435 441 451 546 559 561 580 585 598 612 619 626 635 649 680 701 738 786
[39] 792 795 824 832 834 856 866 868 879 888 907 909 939 949 962 969 971 999

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

[1] 13

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

[1] 9

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

[1] 8

```

## 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/RtmpLXlr/Rinst2fc1449f5f6/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] 49

> sum(diff_results$bootstrap_p<=0.05)

```

```

[1] 64

> 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] 1

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

[1] 5

> 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]
887 cg00862290 0.4364052     0.5404716     0.607868     0.5632595
      exmdata5[, 2] exmdata6[, 2] exmdata7[, 2] exmdata8[, 2]
887      0.5025974     0.4011173     0.566467     0.5455298
      diff_results$ordfit_t[, diff_list_perm]
887                         -3.217939
      diff_results$permutation_p[, diff_list_perm]
887                           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]
259 cg00234961 0.04192170    0.04321576    0.05707140    0.05327565
280 cg00260778 0.64319890    0.60488960    0.56735060    0.53150910
743 cg00717862 0.07999436    0.07873347    0.06089359    0.06171374
911 cg00888479 0.07388961    0.07361080    0.10149800    0.09985076
928 cg00901493 0.03737166    0.03903724    0.04684618    0.04981432
      exmdata5[, 2] exmdata6[, 2] exmdata7[, 2] exmdata8[, 2]
259      0.04030003   0.03996053   0.05086962   0.05445672
280      0.61920530   0.61925200   0.46753250   0.55632410
743      0.07594936   0.09062161   0.06475791   0.07271878
911      0.08633986   0.06765189   0.09070268   0.12417730
928      0.04490690   0.04204062   0.05050039   0.05268215
      diff_results$ordfit_t[, diff_list_boot]

```

```
259          -4.052697
280           4.170347
743            3.444684
911           -3.621731
928           -2.716443
diff_results$bootstrap_p[diff_list_boot]
259             0
280             0
743             0
911             0
928             0
```