# Introduction to RBM package

Dongmei Li

October 17, 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 propor-
  tion 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] 39

> which(myresult$permutation_p<=0.05)

 [1]   22   37   55   56 130 133 163 184 198 227 245 261 298 310 318 320 326 361 474
[20] 476 496 580 600 635 648 698 717 731 757 767 773 802 809 848 930 940 955 961
[39] 977

> sum(myresult$bootstrap_p<=0.05)

[1] 19

> which(myresult$bootstrap_p<=0.05)
```

```
  [1]   10   17   36   37   39   83  144  184  205  276  310  557  719  768  775  809  876  889  940

> 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] 20

> which(myresult2$bootstrap_p<=0.05)

 [1]   24   27   53   79  105  119  175  270  278  287  371  383  494  660  696  734  805  934  976
[20]  980

> 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] 39

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

[1] 44

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

[1] 50

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

 [1]   57   59   89   99  144  191  198  259  261  284  286  295  330  359  365  368  381  402  481
[20]  514  521  560  567  609  632  693  742  743  802  805  807  821  846  896  956  973  979  984
[39]  995

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

 [1]   57   59   99  144  191  198  261  270  284  286  295  320  330  344  359  368  402  473  481
[20]  514  517  521  544  560  567  584  609  632  685  693  742  777  802  805  807  821  846  865
[39]  927  956  973  979  984  995

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

 [1]   49   57   59   99  144  191  198  259  261  270  284  286  295  320  330  343  344  365  368
[20]  381  402  441  473  481  482  514  517  521  546  560  567  605  609  685  693  742  777  802
[39]  805  807  846  865  896  927  947  956  973  979  984  995

> 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] 3

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

[1] 12

> which(con2_adjp<=0.05/3)

[1]   59  284  973
```

```
> which(con3_adjp<=0.05/3)

 [1]  59 191 198 284 286 521 560 609 693 802 973 984

> 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] 48

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

[1] 46

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

[1] 44

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

 [1]   17   57   69   91   98 107 109 113 181 184 194 234 239 241 290 343 351 371 388
[20] 401 419 426 435 453 461 463 464 591 598 600 601 607 609 665 703 745 767 769
[39] 801 834 845 859 864 886 897 909 913 936

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

 [1]   17   19   57   91 107 109 113 181 184 194 239 241 290 312 343 351 371 388 426
[20] 435 453 461 463 464 488 591 598 601 607 609 638 703 745 767 769 776 801 826
[39] 834 864 897 909 913 936 959 963

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

 [1]   17   57   91   98 107 109 113 181 184 239 241 290 343 351 371 388 426 435 453
[20] 461 463 464 488 591 600 601 607 665 703 716 741 767 769 801 826 834 845 859
[39] 864 897 909 913 936 969
```

```
> 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] 4

> 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 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] "/private/tmp/Rtmpx6t6ju/Rinst10a2861f822b5/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] 47

> 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] 0
```

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

[1] 2

> 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)

 [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$ordfit_t[diff_list_perm]
[11] diff_results$permutation_p[diff_list_perm]
<0 rows> (or 0-length row.names)

> 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]
83  cg00072216 0.04505377    0.04598964    0.04000674    0.03231534
743 cg00717862 0.07999436    0.07873347    0.06089359    0.06171374
    exmdata5[, 2] exmdata6[, 2] exmdata7[, 2] exmdata8[, 2]
83     0.04965089    0.04833366    0.03466159    0.04390894
743    0.07594936    0.09062161    0.06475791    0.07271878
    diff_results$ordfit_t[diff_list_boot]
83                               2.514109
743                              3.444684
    diff_results$bootstrap_p[diff_list_boot]
83                                         0
743                                        0
```