

Introduction to RBM package

Dongmei Li

October 27, 2023

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

Contents

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

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 Benjamini-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] 11

> which(myresult$permutation_p<=0.05)

[1] 16 47 262 324 353 399 519 686 835 860 940

> sum(myresult$bootstrap_p<=0.05)

[1] 22

> which(myresult$bootstrap_p<=0.05)

[1] 14 70 86 143 183 190 255 285 302 429 438 481 511 577 640 661 746 811 831
[20] 835 882 953

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

[1] 0

> 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] 17

> which(myresult2$bootstrap_p<=0.05)

[1] 72 98 127 174 179 181 196 371 431 526 633 809 814 846 899 921 978

> bootstrap2_adjp <- p.adjust(myresult2$bootstrap_p, "BH")
> sum(bootstrap2_adjp<=0.05)

[1] 1

```

- 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] 60

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

> sum(myresult_F$permutation_p[, 3]<=0.05)
[1] 72

> which(myresult_F$permutation_p[, 1]<=0.05)
[1] 13 25 34 43 50 108 129 133 136 151 188 203 222 225 255
[16] 283 289 292 304 312 321 327 354 370 399 407 436 446 455 477
[31] 485 487 498 513 527 536 538 546 562 572 580 647 654 678 694
[46] 695 744 748 751 789 830 842 848 881 938 939 952 962 981 1000

> which(myresult_F$permutation_p[, 2]<=0.05)
[1] 13 25 34 43 50 65 81 129 133 136 151 152 188 225 248 255 285 289 292
[20] 310 312 321 327 332 354 370 407 436 446 454 458 477 487 493 498 513 527 529
[39] 536 538 546 580 589 608 638 654 678 694 695 744 748 751 761 781 789 842 848
[58] 881 938 939 952 962 981

> which(myresult_F$permutation_p[, 3]<=0.05)
[1] 13 25 26 34 43 50 65 67 81 129 133 136 143 151 188 198 203 222 225
[20] 272 283 289 292 304 312 321 327 332 346 354 362 370 407 446 454 455 477 487
[39] 498 513 527 531 536 538 546 562 572 580 608 647 654 678 685 694 695 744 748
[58] 761 781 789 830 842 848 870 881 906 938 939 952 959 962 981

> 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] 13

> 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] 13 43 50 136 188 289 354 477 538 546 580 678 952

> which(con3_adjp<=0.05/3)

[1] 13 25 34 50 136 188 203 289 487 538 580 695 952 962

> 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] 54

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

[1] 33

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

```

```

[1] 72 102 105 125 134 152 154 185 269 283 290 378 410 437 468 470 485 495 534
[20] 543 556 638 648 662 664 715 716 731 739 756 773 783 841 880 894 895 906 908
[39] 934 935 938 960 972 973 974

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

[1] 52 55 72 94 102 109 119 120 125 134 152 185 217 233 247 254 269 270 283
[20] 364 378 410 468 470 495 534 543 564 574 638 648 656 662 664 713 715 716 731
[39] 756 773 783 808 816 841 880 894 895 906 934 935 938 960 973 974

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

[1] 72 102 105 125 152 185 247 269 283 290 364 378 410 470 495 513 534 543 648
[20] 656 662 664 715 716 756 783 822 880 894 895 938 960 974

> 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] 3

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

[1] 2

```

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/Rtmp1KcFNJ/Rinstf574126294b2/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] 38

```

```

> sum(diff_results$bootstrap_p<=0.05)
[1] 51

> 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] 0

> 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]
19  cg00016968 0.80628480          NA 0.81440820 0.83623180
259 cg00234961 0.04192170 0.04321576 0.05707140 0.05327565
848 cg00826384 0.05721674 0.05612171 0.06644259 0.06358381
851 cg00830029 0.58362500 0.59397870 0.64739610 0.67269640
               exmdata5[, 2] exmdata6[, 2] exmdata7[, 2] exmdata8[, 2]
19      0.80831380 0.73306440 0.82968340 0.84917800
259     0.04030003 0.03996053 0.05086962 0.05445672
848     0.05230160 0.06119713 0.06542751 0.06240686
851     0.50820240 0.34657470 0.66276570 0.64634510
               diff_results$ordfit_t[, diff_list_perm]
19                  -2.446404
259                 -4.052697
848                 -2.314412
851                 -2.841244
               diff_results$permutation_p[, diff_list_perm]
19                      0
259                      0
848                      0
851                      0

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

```

```
[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$ordfit_t[diff_list_boot]
[11] diff_results$bootstrap_p[diff_list_boot]
<0 rows> (or 0-length row.names)
```