

Introduction to RBM package

Dongmei Li

October 29, 2024

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] 21

> which(myresult$permutation_p<=0.05)
[1] 45 93 112 250 324 405 527 548 627 637 721 769 801 807 809 886 893 910 940
[20] 961 975

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

> which(myresult$bootstrap_p<=0.05)
[1] 36 93 712 717 747 899 987

> 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] 27

> which(myresult2$bootstrap_p<=0.05)
[1] 17 19 38 54 113 138 190 191 221 227 276 373 376 417 468 509 534 590 597
[20] 659 691 719 724 890 909 930 978

> 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] 53

> sum(myresult_F$permutation_p[, 2]<=0.05)
[1] 59

> sum(myresult_F$permutation_p[, 3]<=0.05)
[1] 57

> which(myresult_F$permutation_p[, 1]<=0.05)
[1]   6  15  38  63  90 111 134 153 159 171 182 183 192 193 209 218 263 272 275
[20] 327 329 338 365 404 421 443 447 479 526 548 550 552 577 608 610 642 653 683
[39] 687 695 727 730 738 739 755 791 797 814 818 912 930 931 989

> which(myresult_F$permutation_p[, 2]<=0.05)
[1]   6  15  24  38  63  90 111 134 138 153 171 182 183 192 193 203 209 218 263
[20] 272 275 294 327 329 330 338 365 404 412 421 443 447 479 503 526 550 552 571
[39] 577 608 610 642 643 653 673 676 687 695 712 730 738 791 797 811 818 832 912
[58] 930 931

> which(myresult_F$permutation_p[, 3]<=0.05)
[1]   6  15  38  63  90  94 111 134 138 159 171 182 183 192 193 209 218 263 275
[20] 294 327 329 330 338 404 421 443 447 460 479 526 548 550 577 584 610 625 642
[39] 651 664 676 687 695 712 727 730 738 739 791 797 811 814 818 912 930 931 989

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

[1] 6

```

```

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

[1] 6

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

[1] 2

> which(con2_adjp<=0.05/3)

[1] 183 365 479 526 912 931

> which(con3_adjp<=0.05/3)

[1] 183 192

> 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] 70

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

[1] 60

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

[1] 50

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

[1]   2  77  87 126 141 142 147 160 163 182 190 224 247 259 264 281 285 296 330
[20] 362 365 367 370 376 381 382 385 401 406 446 447 450 451 454 465 476 477 481
[39] 490 495 547 552 553 554 565 601 603 623 639 644 665 682 700 747 762 768 774
[58] 816 823 840 848 853 881 883 894 940 943 948 976 991

```

```

> which(myresult2_F$bootstrap_p[, 2]<=0.05)
[1] 2 36 77 87 101 102 125 141 148 160 182 190 224 247 281 296 330 362 367
[20] 370 376 381 385 401 446 447 450 451 454 465 477 495 514 537 547 553 601 623
[39] 644 665 682 700 747 762 768 774 816 818 823 830 837 840 848 853 881 883 894
[58] 940 976 991

> which(myresult2_F$bootstrap_p[, 3]<=0.05)
[1] 2 77 87 102 126 141 142 148 160 190 224 247 259 281 285 296 330 367 370
[20] 376 381 385 401 406 446 447 450 451 454 465 477 490 495 514 547 553 623 632
[39] 644 665 682 747 762 816 823 830 837 848 881 991

> 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] 7

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

[1] 13

```

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/var/folders/db/4tvgx8jx4z3fm1gzlnlw9rc0000gq/T/RtmpAqeN1P/Rinst166404661eb34/RBM"

> 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] 47

> sum(diff_results$permutation_p<=0.05)
[1] 28

> sum(diff_results$bootstrap_p<=0.05)

```

```

[1] 53

> 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] 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[])
> 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]
95  cg00081975 0.03633894    0.04975194    0.06024723    0.05598723
259  cg00234961 0.04192170    0.04321576    0.05707140    0.05327565
280  cg00260778 0.64319890    0.60488960    0.56735060    0.53150910
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]
95      0.04561792    0.05115624    0.06068253    0.06168212
259     0.04030003    0.03996053    0.05086962    0.05445672

```

```
280    0.61920530    0.61925200    0.46753250    0.55632410
911    0.08633986    0.06765189    0.09070268    0.12417730
928    0.04490690    0.04204062    0.05050039    0.05268215
  diff_results$ordfit_t[diff_list_boot]
95                  -2.654324
259                 -2.833203
280                  4.337628
911                 -3.490240
928                 -1.982308
  diff_results$bootstrap_p[diff_list_boot]
95                      0
259                      0
280                      0
911                      0
928                      0
```