

Introduction to RBM package

Dongmei Li

April 7, 2020

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] 44

> which(myresult$permutation_p<=0.05)

[1] 70 165 175 177 181 211 289 303 425 441 485 487 498 511 516 527 528 582 595
[20] 613 643 648 655 661 668 676 693 717 734 753 760 761 769 776 784 812 863 870
[39] 878 879 891 907 968 996

> sum(myresult$bootstrap_p<=0.05)

[1] 9

> which(myresult$bootstrap_p<=0.05)

[1] 177 342 427 668 676 734 761 779 907

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

[1] 2

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

[1] 1

> 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] 9

> which(myresult2$bootstrap_p<=0.05)

[1] 22 116 418 474 477 626 693 744 867

> 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] 51

> sum(myresult_F$permutation_p[, 3]<=0.05)
[1] 50

> which(myresult_F$permutation_p[, 1]<=0.05)
[1]   6  37  84  93 123 137 173 180 191 201 222 227 254 258 266 302 327 360 364
[20] 438 526 548 550 564 565 590 607 612 614 713 717 732 736 749 755 781 795 798
[39] 806 807 810 822 839 858 874 891 897 914 933 937 965 984 993

> which(myresult_F$permutation_p[, 2]<=0.05)
[1]  37  45  93 137 173 180 196 201 222 254 258 261 266 291 310 327 336 360 364
[20] 438 526 548 550 564 565 590 596 607 612 614 646 692 713 717 732 755 773 781
[39] 795 807 810 822 846 858 867 874 924 937 965 984 994

> which(myresult_F$permutation_p[, 3]<=0.05)
[1]  37  84  93 123 135 173 180 191 254 258 266 291 302 327 332 360 364 414 438
[20] 526 547 548 550 565 596 607 612 614 660 713 717 723 732 773 781 795 798 807
[39] 810 822 839 858 867 872 874 897 937 965 984 993

> 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] 13

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

[1] 8

> which(con2_adjp<=0.05/3)

[1] 266 327 364 550 565 607 612 713 717 795 822 965 984

> which(con3_adjp<=0.05/3)

[1] 93 258 438 550 607 717 822 965

> 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] 75

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

[1] 60

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

[1] 53

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

[1] 15 21 22 27 40 64 74 91 104 111 142 153 158 200 220 233 262 264 265
[20] 278 297 315 316 360 370 377 390 412 413 449 489 491 501 508 513 518 520 531
[39] 537 539 550 558 590 603 626 629 634 653 665 670 679 685 692 702 731 743 753
[58] 756 776 780 797 801 805 838 853 860 868 894 916 923 945 947 962 973 976

```

```

> which(myresult2_F$bootstrap_p[, 2]<=0.05)
[1] 5 15 21 22 27 40 74 91 111 142 153 158 166 220 262 264 265 278 283
[20] 297 315 316 356 360 370 413 489 501 513 518 520 526 531 539 558 603 616 626
[39] 629 653 670 679 692 702 743 756 776 780 801 805 838 868 875 894 916 923 945
[58] 947 973 976

> which(myresult2_F$bootstrap_p[, 3]<=0.05)
[1] 21 22 27 74 91 104 111 141 153 200 220 245 262 264 278 297 315 316 377
[20] 413 449 489 501 513 518 520 526 539 558 590 603 626 653 679 685 692 756 776
[39] 780 801 805 838 860 875 894 916 923 939 945 947 973 976 990

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

[1] 10

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

[1] 6

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

[1] 6

```

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/RtmpKbpWOB/Rinst2188443210a4/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] 42

> sum(diff_results$bootstrap_p<=0.05)

```

```

[1] 62

> 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] 11

> 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]
131 cg00121904 0.15449580    0.17949750    0.23608110    0.24354150
146 cg00134539 0.61101320    0.53321780    0.45999340    0.46787420
200 cg00183916 0.03525946    0.03984548    0.02765822    0.02789838
280 cg00260778 0.64319890    0.60488960    0.56735060    0.53150910
285 cg00263760 0.09050395    0.10197760    0.14801710    0.12242400
518 cg00500400 0.07857063    0.08464774    0.06978949    0.06394599
743 cg00717862 0.07999436    0.07873347    0.06089359    0.06171374
882 cg00858899 0.11427700    0.11919540    0.07690343    0.08321229

```

```

911 cg00888479 0.07388961    0.07361080    0.10149800    0.09985076
928 cg00901493 0.03737166    0.03903724    0.04684618    0.04981432
979 cg00945507 0.13432250    0.23854600    0.34749760    0.28903340
  exmdata5[, 2] exmdata6[, 2] exmdata7[, 2] exmdata8[, 2]
131   0.17352980    0.12564280    0.18193170    0.20847670
146   0.67191510    0.63137380    0.47929610    0.45428300
200   0.03034811    0.04302129    0.02753873    0.03067437
280   0.61920530    0.61925200    0.46753250    0.55632410
285   0.11693600    0.10650430    0.12281160    0.12310430
518   0.07671987    0.07573823    0.07795230    0.06212187
743   0.07594936    0.09062161    0.06475791    0.07271878
882   0.08961409    0.10730660    0.09203980    0.08726349
911   0.08633986    0.06765189    0.09070268    0.12417730
928   0.04490690    0.04204062    0.05050039    0.05268215
979   0.11848510    0.16653850    0.30718420    0.26624740
  diff_results$ordfit_t[diff_list_boot]
131                   -3.451679
146                   5.394750
200                   2.272449
280                   4.170347
285                   -3.093997
518                   2.249342
743                   3.444684
882                   3.179415
911                   -3.621731
928                   -2.716443
979                   -4.750997
  diff_results$bootstrap_p[diff_list_boot]
131                     0
146                     0
200                     0
280                     0
285                     0
518                     0
743                     0
882                     0
911                     0
928                     0
979                     0

```