

Introduction to RBM package

Dongmei Li

May 19, 2021

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] 16

> which(myresult$permutation_p<=0.05)

[1] 27 52 72 190 232 287 291 301 319 499 579 727 773 774 903 915

> sum(myresult$bootstrap_p<=0.05)

[1] 6

> which(myresult$bootstrap_p<=0.05)

[1] 81 454 586 761 895 999

> 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] 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] 11

> which(myresult2$bootstrap_p<=0.05)

[1] 35 67 75 124 317 358 373 429 679 752 830

> 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] 45

> sum(myresult_F$permutation_p[, 2]<=0.05)
[1] 56

> sum(myresult_F$permutation_p[, 3]<=0.05)
[1] 64

> which(myresult_F$permutation_p[, 1]<=0.05)
[1]  2   6 101 105 106 138 145 160 172 178 202 215 224 237 284 288 299 323 367
[20] 408 422 460 479 485 487 499 530 538 542 566 591 614 627 651 675 702 735 761
[39] 780 882 888 903 926 944 988

> which(myresult_F$permutation_p[, 2]<=0.05)
[1]  2   6   28   64   85 101 105 106 129 132 138 145 160 172 202 213 215 237 279
[20] 288 299 323 367 408 420 422 426 460 479 482 485 487 526 530 538 542 566 591
[39] 651 661 702 724 735 761 780 849 882 888 903 924 926 944 945 962 966 988

> which(myresult_F$permutation_p[, 3]<=0.05)
[1]  2   6   78   85 101 105 106 129 138 145 160 172 178 202 203 206 213 215 224
[20] 237 272 279 284 288 299 323 333 367 408 420 422 460 479 485 487 499 526 530
[39] 538 542 566 576 591 627 651 661 675 702 735 761 780 805 832 882 888 903 924
[58] 926 944 945 955 988 993 994

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

[1] 3

```

```

> 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] 7

> which(con2_adjp<=0.05/3)

[1] 106 172 651

> which(con3_adjp<=0.05/3)

[1] 288 323 538 542 735 903 944

> 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] 55

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

[1] 51

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

[1] 51

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

[1]    7   21   22   50   51   60   85  150  181  194  200  205  208  230  244  248  254  255  364
[20] 369  407  412  413  426  446  459  484  487  497  501  513  515  523  526  555  581  603  620
[39] 646  662  669  678  686  703  711  790  832  880  891  907  914  917  941  955  968

```

```

> which(myresult2_F$bootstrap_p[, 2]<=0.05)
[1] 19 21 22 37 46 50 51 85 134 150 181 184 194 200 205 208 244 248 255
[20] 359 364 369 407 412 426 484 487 497 501 513 515 523 526 555 603 620 646 678
[39] 684 703 711 751 832 867 880 907 914 917 926 955 968

> which(myresult2_F$bootstrap_p[, 3]<=0.05)
[1] 7 21 22 37 50 51 85 150 181 194 200 205 208 244 254 255 364 369 407
[20] 412 426 446 459 484 487 497 501 513 523 526 555 556 581 603 618 620 638 669
[39] 678 684 703 711 721 832 880 893 907 914 917 955 968

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

[1] 8

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

[1] 11

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

[1] 10

```

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/RtmpoCBAgF/Rinst2e9527fc7995/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] 64

> 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] 10

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

[1] 4

> 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)

  IlmnID      Beta exmdata2[, 2] exmdata3[, 2] exmdata4[, 2]
19  cg00016968 0.80628480          NA 0.81440820 0.83623180
103 cg00094319 0.73784280 0.73532960 0.75574900 0.73830220
285 cg00263760 0.09050395 0.10197760 0.14801710 0.12242400
627 cg00612467 0.04777553 0.03783457 0.05380982 0.05582291
764 cg00730260 0.90471270 0.90542290 0.91002680 0.91258610
848 cg00826384 0.05721674 0.05612171 0.06644259 0.06358381
851 cg00830029 0.58362500 0.59397870 0.64739610 0.67269640
887 cg00862290 0.43640520 0.54047160 0.60786800 0.56325950
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]
19    0.80831380 0.73306440 0.82968340 0.84917800
103   0.67349260 0.73510200 0.75715920 0.78981220
285   0.11693600 0.10650430 0.12281160 0.12310430
627   0.04740551 0.05332965 0.05775211 0.05579710
764   0.90575890 0.88760470 0.90756300 0.90946790
848   0.05230160 0.06119713 0.06542751 0.06240686
851   0.50820240 0.34657470 0.66276570 0.64634510
887   0.50259740 0.40111730 0.56646700 0.54552980
911   0.08633986 0.06765189 0.09070268 0.12417730
928   0.04490690 0.04204062 0.05050039 0.05268215

  diff_results$ordfit_t[diff_list_perm]
19                      -2.446404
103                     -2.268711

```

```

285           -3.093997
627           -2.239498
764           -1.808081
848           -2.314412
851           -2.841244
887           -3.217939
911           -3.621731
928           -2.716443
diff_results$permutation_p[diff_list_perm]
19              0
103             0
285             0
627             0
764             0
848             0
851             0
887             0
911             0
928             0

> 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]
259 cg00234961 0.04192170   0.04321576   0.05707140   0.05327565
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
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$ordfit_t[diff_list_boot]
259           -4.052697
743            3.444684
911           -3.621731
928           -2.716443
diff_results$bootstrap_p[diff_list_boot]
259              0
743              0
911              0
928              0

```