

Introduction to RBM package

Dongmei Li

May 3, 2016

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:

```
> 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] 29
```

```

> which(myresult$permutation_p<=0.05)
[1] 1 23 71 114 146 180 234 284 306 327 455 487 564 593 605 649 667 688 691
[20] 701 705 708 806 829 851 980 983 985 995

> sum(myresult$bootstrap_p<=0.05)
[1] 0

> which(myresult$bootstrap_p<=0.05)
integer(0)

> 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] 20

> which(myresult2$bootstrap_p<=0.05)
[1] 68 72 152 158 269 308 310 388 470 475 502 696 705 714 741 759 810 869 897
[20] 901

> 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] 66

> sum(myresult_F$permutation_p[, 2]<=0.05)
[1] 55

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

> which(myresult_F$permutation_p[, 1]<=0.05)
[1]   4   5   9  28  49  62  70  78  88 119 123 130 142 145 164 180 184 219 230
[20] 245 246 247 249 267 279 282 290 291 335 346 373 374 375 376 392 393 406 409
[39] 444 461 466 470 473 484 488 505 507 512 580 615 647 671 694 708 771 798 864
[58] 873 878 890 898 927 938 948 951 963

> which(myresult_F$permutation_p[, 2]<=0.05)
[1]   4   5   9  28  29  62  78  88 119 130 145 161 192 218 245 246 249 267 282
[20] 290 291 335 373 374 375 376 393 406 418 433 444 470 484 488 507 512 580 615
[39] 647 651 671 694 708 766 771 798 824 873 878 890 898 927 938 948 963

> which(myresult_F$permutation_p[, 3]<=0.05)
[1]   4   28  62  70  88 119 123 130 145 161 219 245 249 267 279 282 290 335 346
[20] 373 374 375 376 378 392 393 406 412 444 461 470 473 507 512 551 615 647 657
[39] 671 694 708 771 798 873 875 898 927 938 948 963

> 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] 9

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

[1] 10

> which(con2_adjp<=0.05/3)

[1] 28 119 282 373 615 671 694 708 771

> which(con3_adjp<=0.05/3)

[1] 282 373 375 376 392 393 615 771 938 948

> 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] 45

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

[1] 41

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

[1] 36 40 50 55 60 66 74 146 163 183 184 190 191 222 250 255 259 295 305
[20] 319 320 337 360 385 395 413 443 466 468 518 569 580 592 641 645 657 661 675
[39] 688 715 800 824 826 829 832 871 953 972

```

```

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

[1] 1 36 40 53 55 60 66 183 184 191 222 250 255 259 287 295 319 360 395
[20] 413 443 466 468 518 553 558 569 571 592 657 661 675 715 800 806 810 824 829
[39] 847 871 873 953 963 972 980

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

[1] 36 40 50 55 60 66 183 184 191 222 250 255 259 295 314 319 320 337 360
[20] 376 395 413 466 486 518 571 592 599 657 661 669 675 688 715 800 806 824 829
[39] 871 953 972

> 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] 8

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

[1] 4

```

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] "/tmp/RtmpOrRP7P/Rinst3dc16b59a416/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] 58

> sum(diff_results$bootstrap_p<=0.05)

```

```

[1] 44

> 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] 6

> 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)
> 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
627 cg00612467 0.04777553        0.03783457  0.05380982  0.05582291
764 cg00730260 0.90471270        0.90542290  0.91002680  0.91258610
851 cg00830029 0.58362500        0.59397870  0.64739610  0.67269640
887 cg00862290 0.43640520        0.54047160  0.60786800  0.56325950
                           exmdata5[, 2] exmdata6[, 2] exmdata7[, 2] exmdata8[, 2]
19      0.80831380     0.73306440     0.82968340     0.8491780
103     0.67349260     0.73510200     0.75715920     0.7898122
627     0.04740551     0.05332965     0.05775211     0.0557971
764     0.90575890     0.88760470     0.90756300     0.9094679
851     0.50820240     0.34657470     0.66276570     0.6463451
887     0.50259740     0.40111730     0.56646700     0.5455298
    diff_results$ordfit_t[diff_list_perm]
19                  -2.446404
103                 -2.268711
627                 -2.239498
764                 -1.808081
851                 -2.841244
887                 -3.217939
    diff_results$permutation_p[diff_list_perm]
19                      0
103                     0
627                     0

```

```
764          0
851          0
887          0

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