

Introduction to RBM package

Dongmei Li

October 30, 2017

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] 71
```

```

> which(myresult$permutation_p<=0.05)
[1] 2 25 26 28 29 46 70 82 105 118 121 153 156 159 160 176 186 197 202
[20] 203 213 221 222 224 243 294 335 347 353 368 376 393 408 418 438 461 468 578
[39] 581 628 631 635 646 667 717 722 727 732 736 748 750 752 759 760 808 832 839
[58] 841 868 901 903 920 926 932 934 949 970 978 981 991 993

> sum(myresult$bootstrap_p<=0.05)
[1] 11

> which(myresult$bootstrap_p<=0.05)
[1] 82 129 153 160 379 426 567 641 732 736 801

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

[1] 10

> 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] 28

> which(myresult2$bootstrap_p<=0.05)
[1] 70 103 109 148 149 186 197 213 219 233 301 333 428 524 667 697 753 815 827
[20] 831 858 893 907 937 956 960 973 977

> 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] 58

> sum(myresult_F$permutation_p[, 2]<=0.05)
[1] 81

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

> which(myresult_F$permutation_p[, 1]<=0.05)
[1]  33  50  54  58  60  76  78  92  96 135 167 234 236 245 248 284 285 294 312
[20] 332 363 396 401 402 410 412 433 442 449 487 496 499 515 523 535 562 567 575
[39] 579 599 619 626 630 631 675 682 685 697 734 780 821 838 841 861 865 885 938
[58] 980

> which(myresult_F$permutation_p[, 2]<=0.05)
[1]  10  33  39  54  56  58  60  76  78  92 135 167 200 210 223 234 236 245 248
[20] 263 284 285 294 301 312 318 332 360 363 396 402 410 412 427 442 449 463 487
[39] 496 499 519 535 545 547 561 562 567 575 579 599 619 622 626 630 631 668 681
[58] 682 685 697 734 759 780 782 812 821 823 830 838 841 844 845 849 861 865 885
[77] 895 938 939 968 980

> which(myresult_F$permutation_p[, 3]<=0.05)
[1]  33  54  58  60  76  92  96 135 167 210 223 225 234 245 248 284 285 301 312
[20] 313 318 327 332 360 396 397 402 411 412 442 449 487 495 499 535 562 567 575
[39] 579 599 619 626 630 631 657 668 675 681 682 697 734 780 782 812 821 823 830
[58] 838 844 849 861 938 939 980

```

```

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

[1] 12

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

[1] 14

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

[1] 11

> which(con2_adjp<=0.05/3)

[1] 33 58 92 248 294 332 396 410 487 599 630 734 780 938

> which(con3_adjp<=0.05/3)

[1] 54 76 92 248 332 402 442 599 630 682 861

> 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] 47

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

[1] 55

```

```

> which(myresult2_F$bootstrap_p[, 1]<=0.05)
[1] 44 78 82 105 157 170 191 192 251 254 261 307 322 327 331 335 343 348 366
[20] 374 386 428 432 436 438 444 448 490 495 496 503 522 531 552 568 581 584 597
[39] 619 626 627 637 693 701 738 776 855 856 862 868 877 886 902 920 992

> which(myresult2_F$bootstrap_p[, 2]<=0.05)
[1] 81 105 127 157 170 179 191 242 254 261 307 322 327 331 343 348 366 376 386
[20] 428 432 436 440 444 448 495 496 522 552 574 581 584 619 624 626 637 691 701
[39] 738 776 855 862 886 902 937 980 992

> which(myresult2_F$bootstrap_p[, 3]<=0.05)
[1] 44 72 78 81 105 157 170 179 192 242 251 254 261 307 322 327 331 335 343
[20] 348 366 376 386 428 432 436 444 448 453 495 496 514 522 538 552 581 584 597
[39] 619 623 624 626 637 650 701 738 815 837 855 856 862 868 886 902 937

> 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] 5

```

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/Rtmp1KmyLn/Rinst184d60889af2/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] 61

> sum(diff_results$bootstrap_p<=0.05)
[1] 61

> 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] 7

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

[1] 7

> diff_list_perm <- which(perm_adjp<=0.05)
> diff_list_boot <- which(boot_adjp<=0.05)
> sig_results_perm <- cbind(ovarian_cancer_methylation[, diff_results$ordfit_t<=0.05], diff_results$permutation_p<=0.05, diff_results$bootstrap_p<=0.05)
> print(sig_results_perm)

      IlmnID      Beta exmdata2[, 2] exmdata3[, 2] exmdata4[, 2]
83  cg00072216 0.04505377    0.04598964    0.04000674    0.03231534
95  cg00081975 0.03633894    0.04975194    0.06024723    0.05598723
349  cg00332745 0.04703361    0.04634372    0.03676908    0.04518837
804  cg00777121 0.04540701    0.05430304    0.04154242    0.04221162
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
          exmdata5[, 2] exmdata6[, 2] exmdata7[, 2] exmdata8[, 2]
83      0.04965089    0.04833366    0.03466159    0.04390894
95      0.04561792    0.05115624    0.06068253    0.06168212
349      0.04975075    0.05253778    0.04444665    0.03717721
804      0.04911277    0.04872797    0.04261405    0.04474881
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
      diff_results$ordfit_t[diff_list_perm]
83                      2.514109
95                     -3.252063

```

```

349           2.165826
804           1.995220
851          -2.841244
887          -3.217939
911          -3.621731
diff_results$permutation_p[diff_list_perm]
83            0
95            0
349           0
804           0
851           0
887           0
911           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]
95  cg00081975 0.03633894   0.04975194   0.06024723   0.05598723
146  cg00134539 0.61101320   0.53321780   0.45999340   0.46787420
259  cg00234961 0.04192170   0.04321576   0.05707140   0.05327565
280  cg00260778 0.64319890   0.60488960   0.56735060   0.53150910
833  cg00814580 0.09348613   0.09619816   0.12010440   0.11534240
851  cg00830029 0.58362500   0.59397870   0.64739610   0.67269640
979  cg00945507 0.13432250   0.23854600   0.34749760   0.28903340
               exmdata5[, 2] exmdata6[, 2] exmdata7[, 2] exmdata8[, 2]
95    0.04561792   0.05115624   0.06068253   0.06168212
146    0.67191510   0.63137380   0.47929610   0.45428300
259    0.04030003   0.03996053   0.05086962   0.05445672
280    0.61920530   0.61925200   0.46753250   0.55632410
833    0.09577040   0.11598850   0.12860890   0.14111200
851    0.50820240   0.34657470   0.66276570   0.64634510
979    0.11848510   0.16653850   0.30718420   0.26624740
diff_results$ordfit_t[diff_list_boot]
95           -3.252063
146            5.394750
259          -4.052697
280            4.170347
833          -3.428319
851          -2.841244
979          -4.750997
diff_results$bootstrap_p[diff_list_boot]
95            0
146            0
259            0
280            0

```

833	0
851	0
979	0