

Introduction to RBM package

Dongmei Li

April 24, 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] 35
```

```

> which(myresult$permutation_p<=0.05)
[1] 10 19 87 117 170 191 281 301 307 309 332 350 356 408 455 478 527 588 619
[20] 641 676 690 731 751 773 782 787 810 828 876 887 930 937 942 989

> 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] 9

> 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] 42

> which(myresult2$bootstrap_p<=0.05)
[1] 9 22 26 55 56 66 128 135 199 234 250 299 308 340 361 425 447 473 481
[20] 482 517 531 537 548 558 561 571 602 607 624 642 719 756 778 787 833 846 867
[39] 916 930 984 987

> 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] 65

> sum(myresult_F$permutation_p[, 2]<=0.05)
[1] 51

> sum(myresult_F$permutation_p[, 3]<=0.05)
[1] 68

> which(myresult_F$permutation_p[, 1]<=0.05)
[1]   3  18  58  66  83  90 105 123 137 143 150 173 176 186 201 230 235 263 351
[20] 357 404 409 419 433 435 436 446 452 459 486 511 533 556 583 607 614 634 643
[39] 674 676 690 695 699 736 750 761 769 781 790 816 859 865 892 901 902 913 917
[58] 932 949 954 971 972 978 993 999

> which(myresult_F$permutation_p[, 2]<=0.05)
[1]   3  18  29  55  58  65  66  83 105 123 143 173 186 201 215 230 235 351 404
[20] 409 419 435 436 452 459 511 536 556 607 634 643 662 674 690 699 769 781 790
[39] 816 841 855 865 901 913 917 939 949 953 954 972 993

> which(myresult_F$permutation_p[, 3]<=0.05)
[1]   3  18  48  55  58  66  83  90 105 123 150 165 173 176 186 201 209 230 235
[20] 303 305 351 357 374 404 409 419 430 435 436 452 459 486 511 533 536 556 560
[39] 607 619 634 643 674 675 676 690 695 699 736 750 761 769 781 816 861 865 882
[58] 901 913 915 932 949 953 954 972 978 993 999

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

```

```

[1] 14

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

[1] 8

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

[1] 17

> which(con2_adjp<=0.05/3)

[1] 18 186 404 452 459 511 954 993

> which(con3_adjp<=0.05/3)

[1] 3 18 66 186 235 404 419 435 452 459 511 699 736 865 913 954 993

> 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] 74

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

[1] 48

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

[1] 53

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

```

```

[1]   6   9  38  41  44  50  56  67  72  79  90 122 164 165 181 187 200 202 213
[20] 227 231 255 271 290 296 312 314 321 329 337 378 409 410 414 424 442 472 479
[39] 483 490 496 504 508 523 532 550 588 607 616 630 635 651 667 682 689 748 754
[58] 772 777 791 816 839 862 902 913 918 928 933 936 959 972 978 990 991

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

[1]   6   9  38  41  56  72  79 113 164 181 231 280 296 312 314 321 329 337 409
[20] 414 424 472 479 483 496 504 508 588 607 614 616 630 635 667 754 772 777 785
[39] 791 816 839 848 862 918 972 978 990 991

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

[1]   6   9  41  50  56  67  72  79 113 122 165 181 213 231 290 296 312 314 321
[20] 329 337 409 410 424 472 479 483 496 504 508 523 588 607 616 630 635 667 689
[39] 754 772 777 791 816 839 848 862 918 936 945 972 978 990 991

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

[1] 11

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

[1] 9

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

[1] 11

```

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] "C:/Users/biocbuild/bbs-3.5-bioc/tmpdir/RtmpYjcatg/Rinst1c184d7beb4/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] 46

> 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] 2

> 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_results$ordfit_t<=0.05], diff_results$ordfit_t[diff_list_perm])
> print(sig_results_perm)

    IlmnID      Beta exmdata2[, 2] exmdata3[, 2] exmdata4[, 2]
764 cg00730260 0.90471270      0.9054229      0.9100268      0.91258610
911 cg00888479 0.07388961      0.0736108      0.1014980      0.09985076
               exmdata5[, 2] exmdata6[, 2] exmdata7[, 2] exmdata8[, 2]
764      0.90575890      0.88760470      0.90756300      0.9094679
911      0.08633986      0.06765189      0.09070268      0.1241773
               diff_results$ordfit_t[diff_list_perm]
764                           -1.808081
911                           -3.621731
               diff_results$permutation_p[diff_list_perm]
764                               0
911                               0

> sig_results_boot <- cbind(ovarian_cancer_methylation[, diff_results$ordfit_t<=0.05], diff_results$ordfit_t[diff_list_boot])
> print(sig_results_boot)

    IlmnID      Beta exmdata2[, 2] exmdata3[, 2] exmdata4[, 2]
259 cg00234961 0.04192170      0.04321576      0.0570714      0.05327565
280 cg00260778 0.64319890      0.60488960      0.5673506      0.53150910
285 cg00263760 0.09050395      0.10197760      0.1480171      0.12242400
911 cg00888479 0.07388961      0.07361080      0.1014980      0.09985076
               exmdata5[, 2] exmdata6[, 2] exmdata7[, 2] exmdata8[, 2]

```

```
259 0.04030003 0.03996053 0.05086962 0.05445672
280 0.61920530 0.61925200 0.46753250 0.55632410
285 0.11693600 0.10650430 0.12281160 0.12310430
911 0.08633986 0.06765189 0.09070268 0.12417730
  diff_results$ordfit_t[diff_list_boot]
259 -4.052697
280 4.170347
285 -3.093997
911 -3.621731
  diff_results$bootstrap_p[diff_list_boot]
259 0
280 0
285 0
911 0
```