

MLEXPRESSO: differential expression and methylation analysis

Case study using RTCGA data

Aleksandra Dąbrowska, Alicja Gosiewska

Contents

1	Introduction	1
2	Standard Workflow	1
2.1	Methylation	2
2.1.1	BRCA_methylation_chr17 data set	2
2.1.2	Data preparation	2
2.1.3	Testing	2
2.2	Expression	3
2.2.1	BRCA_mRNAseq_chr17 data set	3
2.2.2	Testing	3
2.3	Comparing test results	4
2.4	Visualization	4

1 Introduction

It is considered that the result of increased methylation is decreased gene expression. While, recent studies suggest that the relationship between methylation and expression is more complex than was previously thought.

MLEXPRESSO is an R package for integrative analyses and visualization of gene expression and DNA methylation data.

Key functions of this package are:

- identification of DMR - differentially methylated regions,
- identification of genes with affected expression,
- identification regions with changes in expression and methylation,
- visualization of identified regions.

The joint modeling and visualization of genes expression and methylation improve interpretability of identified signals.

The methodology is supplemented with example applications to The Cancer Genome Atlas data.

2 Standard Workflow

In this vignette we will work with the data sets containing information about gene expression and methylation for patients with breast cancer. We will analyze differences in methylation and expression for patients with different subtypes of BRCA cancer. To run the examples below you should install MLEXPRESSOdata package.(<https://github.com/geneticsMiNing/MLGenSigdata>)

2.1 Methylation

2.1.1 BRCA_methylation_chr17 data set

In this section, we will work with the methylation level data from TCGA database. Package `MLEXPRESSdata` contains `BRCA_methylation_chr17` dataset. This data set contains information about methylation of CpG probes for patients with breast cancer. Rows of this data set correspond to patients, more precisely, to samples taken from patients. First column `SUBTYPE` corresponds to a subtype of BRCA cancer, next columns correspond to CpG probes. Values inside the table indicate the percentage methylation level of CpG probe for specified sample.

```
library(MLEXPRESSdata)
```

```
library(MLEXPRESSdata)
```

```
head(BRCA_methylation_chr17)[1:5,1:4]
```

```
##                SUBTYPE cg00021527 cg00031162 cg00032227
## TCGA-A1-A0SD-01A-11D-A112-05    LumA 0.03781858 0.7910348 0.006391233
## TCGA-A2-A04N-01A-11D-A112-05    LumA 0.01437552 0.7359370 0.008752293
## TCGA-A2-A04P-01A-31D-A032-05    Basal 0.01360124 0.6967802 0.009442039
## TCGA-A2-A04Q-01A-21D-A032-05    Basal 0.01525656 0.5341244 0.014674247
## TCGA-A2-A04T-01A-21D-A032-05    Basal 0.01167384 0.7378100 0.012251559
```

2.1.2 Data preparation

In this analysis we would like to find genes with different methylation. At first we need to use function `aggregate_probes()`, which generates new data frame with CpG probes aggregated to genes. To this aggregation we use, by default, the Illumina Human Methylation data set.

```
BRCA_methylation_gen <- aggregate_probes(BRCA_methylation_chr17)
```

```
head(BRCA_methylation_gen)[1:5,1:4]
```

```
##                AANAT    AARSD1    AATF    AATK
## TCGA-A1-A0SD-01A-11D-A112-05 0.7148533 0.8625816 0.24294092 0.7835302
## TCGA-A2-A04N-01A-11D-A112-05 0.5850106 0.8355825 0.21367129 0.8466190
## TCGA-A2-A04P-01A-31D-A032-05 0.4495537 0.8786166 0.03277413 0.3417919
## TCGA-A2-A04Q-01A-21D-A032-05 0.7120650 0.8819490 0.03460160 0.7264985
## TCGA-A2-A04T-01A-21D-A032-05 0.6010397 0.7739978 0.02501599 0.6276399
```

Before we go to the testing, we need to define condition values for each sample. We would like to test for differences between LumA subtype and other subtypes of breast cancer, so we create a vector, which each element corresponds to a sample. Our division into this two groups relies on numbers of occurrences of each subtype. The LumA subtype is the most common, in case of breast cancer.

```
condition_met <- ifelse(BRCA_methylation_chr17$SUBTYPE=="LumA", "LumA", "other")
```

```
head(condition_met, 8)
```

```
## [1] "LumA" "LumA" "other" "other" "other" "other" "LumA" "other"
```

2.1.3 Testing

In the `MLEXPRESSdata` package we carry out the tests for identification of differentially methylated regions. To do this we use the `calculate_test()` function. Possible values of parameter `test` are described in function documentation.

```
res_met <- calculate_test(BRCA_methylation_gen, condition_met, test="ttest")
head(res_met)
```

```
##      id log2.fold      pval mean_LumA mean_other      mean
## 1  ICAM2 -0.15151320 3.754116e-17 0.2547275 0.4062407 0.3330801
## 2   RILP -0.05073691 2.575168e-13 0.3079069 0.3586438 0.3341447
## 3  PIPOX 0.11505558 5.360053e-12 0.4242804 0.3092248 0.3647812
## 4 TNFSF12 -0.13412855 5.867083e-12 0.1791401 0.3132686 0.2485025
## 5    CD7 0.09822690 1.641919e-11 0.8635077 0.7652808 0.8127112
## 6   KSR1 0.19973400 2.054467e-11 0.658270 0.458536 0.5549808
```

2.2 Expression

2.2.1 BRCA_mRNAseq_chr17 data set

Package `MLEXPRESSdata` contains `BRCA_mRNAseq_chr17` dataset. This set contains information about gene expression: read counts per-gene, computed for genes for 736 patients with breast cancer. Rows of this data set correspond to samples taken from patients. First column `SUBTYPE` corresponds to a subtype of BRCA cancer, next columns correspond to genes.

```
BRCA_mRNAseq_chr17[1:5,1:5]
```

```
##      SUBTYPE AANAT AARSD1 AATF AATK
## TCGA-A1-AOSB-01A-11R-A144-07 Normal      9 2354 2870 317
## TCGA-A1-AOSD-01A-11R-A115-07 LumA      2 1846 5656 312
## TCGA-A1-AOSE-01A-11R-A084-07 LumA     11 3391 9522 736
## TCGA-A1-AOSF-01A-11R-A144-07 LumA      0 2169 4625 169
## TCGA-A1-AOSG-01A-11R-A144-07 LumA      1 2273 3473  92
```

In our example we will test for differential expression between groups with `LumA` breast cancer subtype and `other` subtypes of that cancer. Again we will use vector `conditions`, which consist of two values corresponds to subtype of breast cancer: `LumA` and `other`.

```
condition_exp <- ifelse(BRCA_mRNAseq_chr17$SUBTYPE=="LumA", "LumA", "other")
head(condition_exp, 8)
```

```
## [1] "other" "LumA" "LumA" "LumA" "LumA" "LumA" "other" "LumA"
```

2.2.2 Testing

In the `MLEXPRESS` package we carry out the tests for identification of genes with affected expression. To do this we use the `calculate_test()` function. Possible values of parameter `test` are described in function documentation.

```
res_exp <- calculate_test(BRCA_mRNAseq_chr17[, -1], condition_exp, test="lrt")
head(res_exp)
```

```
##      id log2.fold      pval mean_LumA mean_other      mean
## 1  AURKB 2.339920 3.191000e-32 539.0426 2323.8868 1485.01
## 2   CBX2 2.895062 2.834335e-26 632.5106 4296.6038 2574.48
## 3  KPNA2 1.447288 8.551812e-24 11547.36 26427.38 19433.77
## 4  PRR11 3.822148 2.286874e-22 396.383 3479.981 2030.69
## 5  BIRC5 1.988998 1.953941e-21 1957.085 6658.358 4448.76
## 6   GSG2 1.405039 3.527773e-21 278.2128 629.3396 464.31
```

2.3 Comparing test results

We can also create a comparison table with results of `calculate_test()` function for methylation and expression data. With this two results we compute the ranking of the most significant changed genes in terms of both methylation and expression. The created column contains the geometric mean of p-values for expression and methylation.

```
genes_comparison <- calculate_comparison_table(BRCA_mRNAseq_chr17[, -1], BRCA_methylation_gen, condition)

## Warning in sqrt(result[, 2] * result[, 4]): NaNs produced
## Warning: Column `id` joining character vector and factor, coercing into
## character vector
head(genes_comparison)

##           id nbinom2.log2.fold nbinom2.pval ttest.log2.fold ttest.pval
## 354    LSM12      0.056334003    0.4954928    9.234253e-05 0.87094671
## 579  SMARCE1     -0.022895312    0.8318716   -2.377235e-04 0.60251281
## 77   C17orf61    -0.006474155    0.9585554   -1.027549e-03 0.10515489
## 367    MED9     -0.003410034    0.9677834   -2.615948e-03 0.03653572
## 482  PRKAR1A    -0.003362801    0.9791396   -4.099641e-03 0.08891401
## 273    GRB2     -0.067708427    0.5412275   -2.148793e-04 0.85842993
##      geom.mean.rank No.probes
## 354    0.002280795      2
## 579    0.002332971      2
## 77     0.002579247      2
## 367    0.002986716      1
## 482    0.003712987      2
## 273    0.003814334      2
```

2.4 Visualization

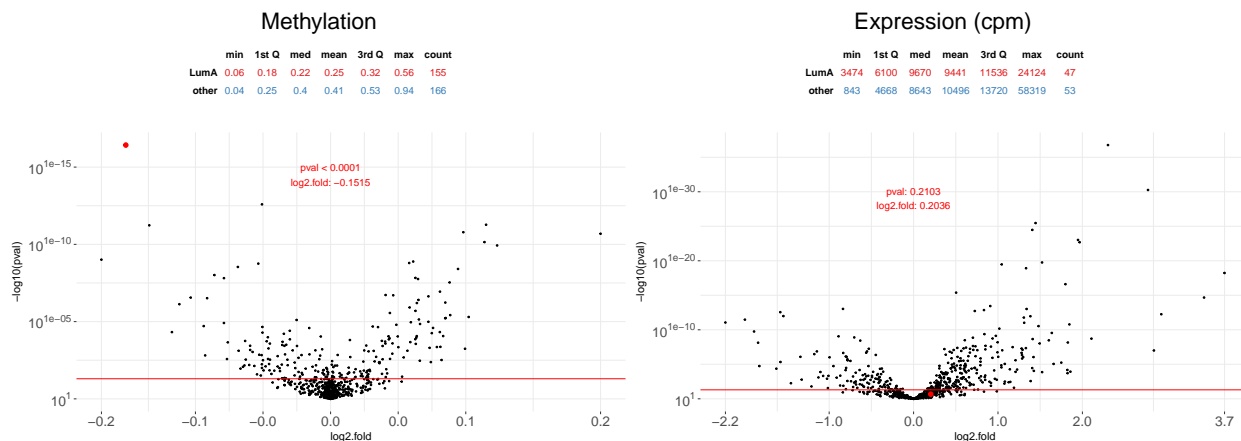
```
test_exp <- genes_comparison[, c(1, 2, 3)]
test_met <- genes_comparison[, c(1, 4, 5)]
```

The great advantage of MLEXPRESSO package is the ability to perform a variety of visualizations for expression and methylation.

For both, methylation and expression data, we can visualise the volcano plots for results of chosen tests and simple statistics for chosen gene.

```
plot_volcanoes(BRCA_methylation_chr17[, -1], BRCA_mRNAseq_chr17[, -1], condition_met, condition_exp,
               "ICAM2",
               test_met, test_exp,
               values=TRUE)
```

ICAM2



Other function `plot_gene()` allow us to visualise the **methylation path** - placement of probes near the gene with a marked percentage of methylation for each probe in division into groups. Using this function we also get boxplots containing values from expression in division from `condition_exp` vector for chosen gene. Note that `plot_gene()` methylation require data frame with CpG probes, not genes.

```
plot_gene(BRCA_methylation_chr17, BRCA_mRNAseq_chr17, condition_met, condition_exp, "ICAM2")
```

