

Using MBatch Corrections: MP_ByBatch

Tod Casasent

2018-05-03

1 Introduction

These instructions are aimed at people familiar with R and familiar with TCGA/GDC platforms and data types. They are intended to introduce the reader to producing the given assessment. These instructions will only rarely, if ever, touch on the appropriateness of the assessment algorithm or interpretation of output. See MBatch_01_InstallLinux for instructions on downloading test data.

2 Algorithm

MP Overall performs a Median Polish Overall correction taking a BEA_DATA object (with data matrix and batch dataframe) and returning either a corrected matrix or a string containing the path to where the data file was written.

3 Output

The primary output method for MBatch is to view results in the Batch Effects Website. Correction algorithms generally do not create graphical output and instead create TSV output files.

4 Usage

```
MP_ByBatch(theBeaData, theBatchType, thePath = NULL, theWriteToFile = FALSE)
```

5 Arguments

5.1 theBeaData

BEA_DATA objects can be created by calls of the form `new("BEA_DATA", theData, theBatches, theCovariates)`. If you have no covariate data, use an empty data.frame created with `data.frame()`

mData: Object of class "matrix" A matrix where the colnames are sample ids and the rownames are gene equivalents. All names should be strings, not factors.

mBatches: Object of class "data.frame" A data.frame where the column "names" are batch types. The first batch "type" is "Sample". All names and values should be strings, not factors or numeric.

mCovariates: Object of class "data.frame" A data.frame where the column "names" are covariate types. The first covariate "type" is "Sample". All names and values should be strings, not factors or numeric.

5.2 theBatchType

A string identifying the batch type to correct.

5.3 thePath

Output path for any files.

5.4 theWriteToFile

TRUE to write the corrected data to file and return the file pathname instead of the corrected matrix.

6 Example Call

The following code is adapted from the tests/MP_ByBatch.R file. Data used is from the testing data as per the MBatch_01_InstallLinux document. In the future, we plan to make the output from MBatch more user friendly, but currently, this produces the following output at the command line.

```
{
  library(MBatch)

  # set the paths
  invariantFile="/bea_testing/MATRIX_DATA/rbn-pseudo-iset.tsv"
  variantFile="/bea_testing/MATRIX_DATA/rbn-pseudo-vset.tsv"
  theOutputDir="/bea_testing/output/RBN_Pseudoreplicates"
  theRandomSeed=314

  theGeneFile="/bea_testing/MATRIX_DATA/matrix_data-Tumor.tsv"
  theBatchFile="/bea_testing/MATRIX_DATA/batches-Tumor.tsv"
  theOutputDir="/bea_testing/output/MP_ByBatch"
  theRandomSeed=314
  theBatchType="TSS"

  # make sure the output dir exists and is empty
  unlink(theOutputDir, recursive=TRUE)
  dir.create(theOutputDir, showWarnings=FALSE, recursive=TRUE)
  # load data
  myData <- mbatchLoadFiles(theGeneFile, theBatchFile)
  myData@mData <- mbatchTrimData(myData@mData, 100000)
  # call
  outputFile <- MP_ByBatch(theBeaData=myData,
                           theBatchType=theBatchType,
                           thePath=theOutputDir,
                           theWriteToFile=TRUE)
  correctedMatrix <- readAsGenericMatrix(outputFile)
  print(correctedMatrix[1:4, 1:4])
}
```

```
## 2018 05 03 13:40:51.298 DEBUG megazone23 Changing LC_COLLATE to C for duration of run
## 2018 05 03 13:40:51.299 INFO megazone23 \ / \ / \ / \ / \ / \ / \ / \ / \ /
## 2018 05 03 13:40:51.299 INFO megazone23 Starting mbatchLoadFiles
## 2018 05 03 13:40:51.299 INFO megazone23 MBatch Version: 2017-09-19-1530
## 2018 05 03 13:40:51.299 INFO megazone23 read batch file= /bea_testing/MATRIX_DATA/batches-Tumor.tsv
## 2018 05 03 13:40:51.301 INFO megazone23 read gene file= /bea_testing/MATRIX_DATA/matrix_data-Tumor.t
## 2018 05 03 13:41:03.204 INFO megazone23 filter samples in batches using gene samples
## 2018 05 03 13:41:03.205 INFO megazone23 sort batches by gene file samples
```

```

## 2018 05 03 13:41:03.301 INFO megazone23 Finishing mbatchLoadFiles
## 2018 05 03 13:41:03.301 INFO megazone23 ~~~~~
## 2018 05 03 13:41:03.301 DEBUG megazone23 Changing LC_COLLATE to C for duration of run
## 2018 05 03 13:41:03.301 INFO megazone23 \ / \ / \ / \ / \ / \ / \ / \ / \ /
## 2018 05 03 13:41:03.302 INFO megazone23 mbatchTrimData Starting
## 2018 05 03 13:41:03.302 INFO megazone23 MBatch Version: 2017-09-19-1530
## 2018 05 03 13:41:10.485 INFO megazone23 mbatchTrimData Finishing
## 2018 05 03 13:41:10.485 INFO megazone23 ~~~~~
## 2018 05 03 13:41:10.485 INFO megazone23 MP_Internal - starting
## 2018 05 03 13:41:10.486 DEBUG megazone23 checkCreateDir: /bea_testing/output/MP_ByBatch
## 2018 05 03 13:41:10.655 DEBUG megazone23 starting BeaMP
## 2018 05 03 13:41:10.655 DEBUG megazone23 starting MP
## 2018 05 03 13:41:10.656 DEBUG megazone23 MP batch
## 2018 05 03 13:41:10.656 DEBUG megazone23 convertDataFrameToSi start
## 2018 05 03 13:41:10.656 DEBUG megazone23 convertDataFrameToSi asmatrixWithIssues
## 2018 05 03 13:41:10.658 DEBUG megazone23 convertDataFrameToSi rownames
## 2018 05 03 13:41:10.659 DEBUG megazone23 convertDataFrameToSi colnames
## 2018 05 03 13:41:10.659 DEBUG megazone23 convertDataFrameToSi done
## 2018 05 03 13:41:12.020 DEBUG megazone23 finishing BeaMP
## 2018 05 03 13:41:12.020 TIMING megazone23      1.3679999999999999      1.3659999999999999      MPByBatch /b
## 2018 05 03 13:41:12.021 DEBUG megazone23 Write to file /bea_testing/output/MP_ByBatch/ANY_Correction
## 2018 05 03 13:41:12.123 DEBUG megazone23 Finished write to file /bea_testing/output/MP_ByBatch/ANY_
## 2018 05 03 13:41:12.123 INFO megazone23 MP_Internal - completed
##
##          TCGA-OR-A5J1-01A-11D-A29J-05
## ABR-cg06968724-17-1012579      -0.3519358
## ABR-cg23568341-17-1011974      -0.3830771
## ABR-cg24479027-17-1012576      -0.3482694
## ACOT7-cg16034168-1-6336711      0.4385622
##
##          TCGA-OR-A5J2-01A-11D-A29J-05
## ABR-cg06968724-17-1012579      0.4393437
## ABR-cg23568341-17-1011974      0.4134640
## ABR-cg24479027-17-1012576      0.4451167
## ACOT7-cg16034168-1-6336711      0.3473410
##
##          TCGA-OR-A5J3-01A-11D-A29J-05
## ABR-cg06968724-17-1012579      0.97342630
## ABR-cg23568341-17-1011974      0.87449504
## ABR-cg24479027-17-1012576      0.98826342
## ACOT7-cg16034168-1-6336711      0.01531205
##
##          TCGA-OR-A5J4-01A-11D-A29J-05
## ABR-cg06968724-17-1012579      0.5599403
## ABR-cg23568341-17-1011974      0.5152550
## ABR-cg24479027-17-1012576      0.5571905
## ACOT7-cg16034168-1-6336711      0.4181180

```

7 Example File Output

The above code creates the following output file. File is named using the following naming convention: ANY_Corrections-MPByBatch.tsv The TSV file with the corrected dataset is written by the MBatch package. The end of the output shows a snippet from the corrected matrix.