

## 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.docx for instructions on downloading test data.

## 2 Algorithm

EBNPlus\_CombineBatches is a function used to combine batch information after the data has been combined via the EBNPlus algorithm.

## 3 Output

The primary output method for MBatch is to view results in the Batch Effects Website. Correction algorithms and the combine batches function does not create graphical output and instead creates TSV output files.

## 4 Usage

```
CDP_Structures(theFilePath, theData1, theData2, theSubTitle,  
               theUnmatchedCount = 1000, theMethod = "pearson",  
               theUse = "pairwise.complete.obs", theSeed = NULL,  
               theUseReplicatesUnpaired=FALSE,  
               theLinePlot=TRUE, theHistPlot=TRUE, theBinWidth=NULL)
```

## 5 Arguments

**theFilePath** Full path and filename for PNG output file

**theData1** Matrix with columns as samples and rows as features.

**theData2** Matrix with columns as samples and rows as features.

**theSubTitle** Subtitle for image, giving data type being displayed.

**theUnmatchedCount** Number of iterations for unpaired samples.

**theMethod** Defaults to pearson. Valid values are: concordance, pearson, kendall, spearman.

**theUse** Defaults to pairwise.complete.obs. Valid values are accepted by the method parameter to cor.

**theSeed** Default to NULL.

**theUseReplicatesUnpaired** Defaults to FALSE. If TRUE, use both the replicates and non-replicates for the unpaired plot.

**theLinePlot** Default to TRUE. TRUE means plot the lines for Correlation Density Plots.

**theHistPlot** Default to TRUE. TRUE means plot the histogram for Correlation Density Plots.

**theBinWidth** Default to NULL. Non-null means to use the given wide for bins. Otherwise, use default for hist.

## 6 Example Call

The following code combined batch files and is taken from the tests/EBNPlus\_CombineBatches.R file. Data used is from the testing data as per the MBatch\_01\_InstallLinux.docx document.

```
library(MBatch)

# set the paths
theBatchFile="/bea_testing/MATRIX_DATA/brca_rnaseq2_batches.tsv"
theBatchFile2="/bea_testing/MATRIX_DATA/brca_agi4502_batches.tsv"
theOutputDir="/bea_testing/output/EBNPlus_CombineBatches"
theBatchId1="RNASeqV2"
theBatchId2="Agilent4502"

# make sure the output dir exists and is empty
unlink(theOutputDir, recursive=TRUE)
dir.create(theOutputDir, showWarnings=FALSE, recursive=TRUE)

dataBatches <- EBNPlus_CombineBatches(readAsDataFrame(theBatchFile),
                                       readAsDataFrame(theBatchFile2), theBatchId1, theBatchId2)
writeAsDataframe(file.path(theOutputDir, "BatchData.tsv"), dataBatches)
```

### 6.1 Command Line Output

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
> theBatchFile="/bea_testing/MATRIX_DATA/brca_rnaseq2_batches.tsv"
> theBatchFile2="/bea_testing/MATRIX_DATA/brca_agi4502_batches.tsv"
> theOutputDir="/bea_testing/output/EBNPlus_CombineBatches"
> theBatchId1="RNASeqV2"
> theBatchId2="Agilent4502"
>
> # make sure the output dir exists and is empty
> unlink(theOutputDir, recursive=TRUE)
> dir.create(theOutputDir, showWarnings=FALSE, recursive=TRUE)
>
> dataBatches <- EBNPlus_CombineBatches(readAsDataFrame(theBatchFile),
+                                       readAsDataFrame(theBatchFile2), theBatchId1, theBatchId2)
2017 10 17 13:42:00.570 DEBUG MachineName starting BeaEBNPlusBatches
2017 10 17 13:42:00.572 DEBUG MachineName readAsDataFrame - thePar -Xmx2000m
2017 10 17 13:42:00.572 DEBUG MachineName readAsDataFrame - theFile
/bea_testing/MATRIX_DATA/brca_rnaseq2_batches.tsv
```

```

2017 10 17 13:42:00.572 DEBUG MachineName readAsDataFrame - Calling .jinit
/home/linux/R/x86_64-pc-linux-gnu-library/3.4/MBatch/ReadRJava/ReadRJava.jar
2017 10 17 13:42:00.579 DEBUG MachineName readAsDataFrame - .jinit complete
2017 10 17 13:42:00.580 DEBUG MachineName readAsDataFrame before java
ReadRJavaL::loadStringData 2014-04-20-1523
2017 10 17 13:42:00.682 DEBUG MachineName readAsDataFrame after java
2017 10 17 13:42:00.685 DEBUG MachineName readAsDataFrame - length(myData) 7290
ReadRJavaL::loadStringData done
2017 10 17 13:42:00.686 DEBUG MachineName readAsDataFrame - length(myCols) 6
2017 10 17 13:42:00.686 DEBUG MachineName readAsDataFrame - length(myRows) 0
2017 10 17 13:42:00.687 DEBUG MachineName readAsDataFrame - myCols Sample, Type, BatchId,
PlateId, ShipDate, TSS
2017 10 17 13:42:00.688 DEBUG MachineName readAsDataFrame - myRows
2017 10 17 13:42:00.690 DEBUG MachineName readAsDataFrame - thePar -Xmx2000m
2017 10 17 13:42:00.691 DEBUG MachineName readAsDataFrame - theFile
/bea_testing/MATRIX_DATA/brca_agi4502_batches.tsv
2017 10 17 13:42:00.691 DEBUG MachineName readAsDataFrame - Calling .jinit
/home/linux/R/x86_64-pc-linux-gnu-library/3.4/MBatch/ReadRJava/ReadRJava.jar
2017 10 17 13:42:00.698 DEBUG MachineName readAsDataFrame - .jinit complete
2017 10 17 13:42:00.699 DEBUG MachineName readAsDataFrame before java
ReadRJavaL::loadStringData 2014-04-20-1523
ReadRJavaL::loadStringData done
2017 10 17 13:42:00.703 DEBUG MachineName readAsDataFrame after java
2017 10 17 13:42:00.705 DEBUG MachineName readAsDataFrame - length(myData) 3600
2017 10 17 13:42:00.706 DEBUG MachineName readAsDataFrame - length(myCols) 6
2017 10 17 13:42:00.706 DEBUG MachineName readAsDataFrame - length(myRows) 0
2017 10 17 13:42:00.707 DEBUG MachineName readAsDataFrame - myCols Sample, Type, BatchId,
PlateId, ShipDate, TSS
2017 10 17 13:42:00.708 DEBUG MachineName readAsDataFrame - myRows
> writeAsDataframe(file.path(theOutputDir, "BatchData.tsv"), dataBatches)
2017 10 17 13:42:00.747 DEBUG MachineName writeAsDataframe - thePar -Xmx2000m
2017 10 17 13:42:00.748 DEBUG MachineName writeAsDataframe - theFile
/bea_testing/output/EBNPlus_CombineBatches/BatchData.tsv
2017 10 17 13:42:00.749 DEBUG MachineName writeAsDataframe - length(myData) 12705
2017 10 17 13:42:00.749 DEBUG MachineName writeAsDataframe - length(myCols) 7
2017 10 17 13:42:00.750 DEBUG MachineName writeAsDataframe - length(myRows) 0
2017 10 17 13:42:00.750 DEBUG MachineName writeAsDataframe - Calling .jinit
/home/linux/R/x86_64-pc-linux-gnu-library/3.4/MBatch/ReadRJava/ReadRJava.jar
2017 10 17 13:42:00.756 DEBUG MachineName writeAsDataframe - .jinit complete
2017 10 17 13:42:00.757 DEBUG MachineName writeAsDataframe before java
ReadRJava::writeStringData_Column 2014-04-20-1523
writeFile - start
writeFile - done
ReadRJava::writeStringData_Column done
2017 10 17 13:42:00.767 DEBUG MachineName writeAsDataframe after java
2017 10 17 13:42:00.767 DEBUG MachineName writeAsDataframe success= TRUE
[1] TRUE

```

## 6.2 Example File Output

The above code creates the following output files. Files are named using the following naming convention:

BatchData.tsv

This is a TSV file with both original batch files combined and save here.